** VARYING PROBABILITY SAMPALING **

Auxiliary Information:-

- (i) Planning stage e.g. Stratified Sampling
- (ii) Specified Stage
- (iii) Estimation Stage e.g. Ratio, Regression, Product method of estimation.

Unlike SRS, here the units gets different probability of selection with the help of the use auxiliary variable i.e. We associate to vector $Y' = (y_1, \dots, y_n)$, a vector $P' = (p_1, p_2, \dots, p_n)$ such that $\sum_{i=1}^{N} p_i = 1$, $p_i > 0$. The rational behind the use of varying probability of selection is that large units contribute more towards the population total and

probability of selection is that large units contribute more towards the population total and hence they should have large probability of selection in the sample. If X_i denotes the size of i^{th} unit. Then $p_i \propto x_i$

$$\therefore \sum_{i=1}^{N} p_i = k \sum_{i=1}^{N} x_i \implies k = \frac{1}{X} \text{ where } X = \sum_{i=1}^{N} x_i$$
$$\therefore p_i = \frac{x_i}{X} \text{ where } i = 1, 2, \dots, N.$$

This sampling is also known as PPS (Probability Proportional to Size)

There are predominately two methods by which we select the PPS sample.

- (i) Cumulative total method.
- (ii) Laheri`s method.
- (1) Cumulative Total Method:-

Here we make the cumulative total of the units of the population given as bellow,

Unit No.	Size	Cumulative total
1	X1	$T_1 = x_1$
2	X2	$T_2 = x_1 + x_2$
3	X3	$T_3 = x_1 + x_2 + x_3$
•	•	•••••
•	•	
Ι	Xi	$T_i = x_1 {+} {\ldots} {+} x_i$
•		•••
Ν	X _N	$T_N=x_1{+}{\ldots}{+}x_N$

Now select a random number between 1 to X. Where X = total of X_i's in population. If rth random number is selected and $T_{i-1} < r \le T_i$. Then ith unit is selected. This ensure that $p_i = x_i / X$. If we allow repetition of units in the sample. It is known as PPSWR sampling procedure and if we do not allow repetition then resultant will be PPSWOR sampling procedure.

One of the major defect of this sampling procedure is that for large N. Preparation of cumulative total table requires large cost and time, so Laheri suggested the following procedure.

(2) Laheri's method of selection of a PPS sample:-

Laheri's Scheme consists in 1st deciding a number in which it is larger than the sizes of units. Then select a pair of random number say (i, j), if $x_i \le x_j$ retain ith unit in the sample otherwise reject the pair. Repeat it till we get the sample size n. Here also if retain all units even with repetition, it results an PPSWR. If we retain only distinct units in sample, then it is PPSWOR.

** Estimation Of population mean or total.

Let a sample of $(y_1,...,y_n)$ is selected with PPSWR in which probability of selecting i^{th} unit of y is p_i (i = 1, 2, ..., N) and we are interested in population mean.

$$\overline{Y} = \frac{\sum_{i=1}^{N} Y_i}{N}$$

let $Z_i = \frac{Y_i}{NP_i}$ $i = 1, 2, \dots, N$
 $\overline{z} = \frac{1}{n} \sum_{i=1}^{n} z_i = \frac{1}{n} \sum_{i=1}^{n} \frac{y_i}{NP_i}$

then

is an unbiased estimator of population mean, that is,

$$E(\overline{z}) = \frac{1}{n} \sum_{i=1}^{n} E(z_i) = \frac{1}{n} \sum_{i=1}^{n} [\Sigma Z_i p_i] = \frac{1}{n} \sum_{i=1}^{n} \Sigma \frac{Y_i}{NP_i} p_i$$
$$= \frac{1}{n} \sum_{i=1}^{n} \overline{Y} = \frac{1}{n} n \overline{Y} = \overline{Y}$$
$$E(z_i) = \sum_{i=1}^{N} P_i Z_i = \sum_{i=1}^{N} P_i \frac{Y_i}{NP_i} = \sum_{i=1}^{N} Y_i = \overline{Y}$$
$$V(z_i) = E(z_i - \overline{Y})^2 = \sum_{i=1}^{N} P_i (Z_i - \overline{Y})^2 = \sum_{i=1}^{N} P_i (\frac{Y_i}{NP_i} - \overline{Y})^2 = \sigma_z^2, \text{ say}$$

$$V(\overline{z}) = E(\overline{z}^{2}) - \left[E(\overline{z})\right]^{2}$$
(1)

$$E(\overline{z}^{2}) = E\left[\frac{1}{n}\sum_{i=1}^{n} z_{i}\right]^{2}$$
$$= \frac{1}{n^{2}}E\left[\sum_{i=1}^{n} z_{i}\right]^{2}$$
$$= \frac{1}{n^{2}}E\left[\sum_{i} z_{i}^{2} + \sum_{i\neq j}^{n} z_{i}z_{j}\right]$$
$$= \frac{1}{n^{2}}\left[\sum_{i} E(z_{i}^{2}) + \sum_{i\neq j}^{n} E(z_{i})E(z_{j})\right]$$
(ii)

as z_i and z_j are independent.

$$E(z_i^2) = \sum_i z_i^2 p_i$$

$$E(z_i) = \overline{z}$$
 (iii)

On putting this value in (ii), we have

$$E(z^{-2}) = \frac{1}{n^2} \left[\sum_{i} \sum_{i} z_i^2 p_i + \sum_{i \neq j}^n z^{-2} \right]$$
$$= \frac{1}{n^2} \left[n \sum_{i} z_i^2 p_i + n(n-1) z^{-2} \right]$$
$$= \frac{1}{n} \left[\sum_{i} z_i^2 p_i + (n-1) z^{-2} \right]$$
$$= \frac{1}{n} \left[\sum_{i} z_i^2 p_i + nz^{-2} - z^{-2} \right]$$

On putting this value in (i), we have

$$\operatorname{Var}\left(\overline{z}\right) = \frac{1}{n} \left[\sum_{i} z_{i}^{2} p_{i} + nz^{-2} - z^{-2} \right] - z^{-2}$$
$$= \frac{1}{n} \left[\sum_{i} z_{i}^{2} p_{i} + nz^{-2} - z^{-2} - nz^{-2} \right]$$
$$= \frac{1}{n} \left[\sum_{i} z_{i}^{2} p_{i} - z^{-2} \right]$$
$$= \frac{1}{n} \sum_{i} p_{i} (z_{i} - z^{-})^{2}$$
$$= \frac{\sigma^{2} z}{n} \quad \text{where } \sigma^{2} z = \sum_{i} p_{i} (z_{i} - z^{-})^{2}$$

** Unbiased Estimator V($\overline{\chi}$) :-

Let us consider
$$s_z^2 = \frac{1}{n-1} \sum_{1}^{n} \left[(z_i - z^{-})^2 \right]$$

Then, $E(s_z^2) = \frac{1}{n-1} E\left[\sum_{1}^{n} (z_i - z^{-})^2 \right]$
 $(n-1) E(s^2z) = E\left[\sum_{1}^{n} (z_i^2 - nz^{-2}) \right]$
 $= \sum_{1}^{n} E(z_i^2) - nE(z^{-2})$
 $= n \sum_{1} z_i^2 p_i - n \left[Var(z^{-}) + \left[E(z^{-}) \right]^2 \right]$
 $= n \sum_{1} z_i^2 p_i - n \left[\frac{\sigma_z^2}{n} z^{-2} \right]$

$$= n \sum_{i} z_{i}^{2} p_{i} - \sigma_{z}^{2} - n \frac{1}{z^{2}}$$
$$= n \left[\sum_{i} z_{i}^{2} p_{i} - \frac{1}{z^{2}} \right] - \sigma_{z}^{2}$$
$$= n \sigma_{z}^{2} - \sigma_{z}^{2}$$
$$= (n - 1) \sigma_{z}^{2}$$
$$E(s_{z}^{2}) = \sigma_{z}^{2}$$
$$V(\overline{z}) = \sigma_{z}^{2} / n = E(s_{z}^{2}) / n$$

Unbiased estimator of $V(z) = s_z^2/n$

...

** PPSWOR :-

This was first considered Horvits & Thompson (1952). The use of PPWOR for estimating population mean. The suggested that

$$\overline{Y}_{HT} = \sum_{i=1}^{N} \frac{y_i}{N\pi_i}$$

where π_i = inclusion probability of ith unit being included in the sample. This is different from initial probability of selection which is given by p_i. Then,

Where π_{ij} is the probability of including units i and j in the sample.

$$\begin{aligned} \pi_{ij} &= \sum_{s > (i,j)} p(s) \\ v(\overline{Y_{HT}}) &= v \Biggl[\sum_{i=1}^{N} \frac{Y_i s_i}{N \pi_i} \Biggr] \\ &= \frac{1}{N^2} \Biggl[\sum_{i=1}^{N} \frac{Y_i^2 v(s_i)}{\pi_i^2} + \sum_{i \neq j}^{N} \frac{Y_i Y_j \operatorname{cov}(\delta_i, \delta_j)}{\pi_i \pi_j} \Biggr] \\ &= \frac{1}{N^2} \Biggl[\sum_{i=1}^{N} \frac{Y_i^2 \pi_i (1 - \pi_i)}{\pi_i^2} + \sum_{i \neq j}^{N} \frac{Y_i Y_j (\pi_{ij} - \pi_i \pi_j)}{\pi_i \pi_j} \Biggr] \\ &= \frac{1}{N^2} \Biggl[\sum_{i=1}^{N} \frac{Y_i^2 (1 - \pi_i)}{\pi_i} + \sum_{i \neq j}^{N} \frac{Y_i Y_j (\pi_{ij} - \pi_i \pi_j)}{\pi_i \pi_j} \Biggr] \end{aligned}$$

** Estimation of $v(\overline{Y_{HT}})$:-

$$v(\overline{Y_{HT}}) = \frac{1}{N^2} \left[\sum_{i=1}^{N} \frac{\pi_i Y_i^2 (1 - \pi_i)}{\pi_i^2} + \sum_{i \neq j}^{N} \frac{Y_i Y_j (\pi_{ij} - \pi_i \pi_j)}{\pi_i \pi_j} \right]$$
$$= \frac{1}{N^2} \left[\sum_{i=1}^{N} \frac{\pi_i (1 - \pi_i) Y_i^2}{\pi_i^2} + \sum_{i \neq j}^{N} \frac{Y_i Y_j Y_i (\pi_{ij} - \pi_i \pi_j)}{\pi_i \pi_j} \right]$$

* Evaluation of inclusion probability n =2.

$$\begin{split} \pi_{i} &= \sum_{r=1}^{n} p_{ir} \\ &= \sum_{r=1}^{2} p_{ir} \\ &= p_{i1} + p_{i2} \\ &= p_{i} + \sum_{j(\neq i)=1}^{N} \frac{p_{j} p_{i}}{1 - p_{j}} \\ &= p_{i} \left(1 + s - \frac{p_{i}}{1 - p_{i}} \right) \quad \text{where} \quad s = \sum_{i=1}^{N} \frac{p_{i}}{1 - p_{i}} \\ \pi_{ij} &= p_{i1j2} + p_{j1i2} \\ &= \frac{p_{i} p_{j}}{1 - p_{i}} + \frac{p_{j} p_{i}}{1 - p_{j}} \\ &= p_{i} p_{j} \left[\frac{1}{1 - p_{i}} + \frac{1}{1 - p_{j}} \right] \end{split}$$

** Yates and Grondy proved the following relationship.

An elegent expression for the variance of the H-T estimator is given by Yates and Grondy. Since, $\delta_I = 1$ exactly n units in the probability and zero for the rest, we have,

Take expectation of both sides, we have

$$\sum_{i=1}^{N} E(\delta_i) = E(n)$$
$$\sum_{i=1}^{N} \pi_i = n$$

Squaring (1) and taking the expectation, we have

$$\sum_{i\neq j}^{N} \pi_{ij} = n(n-1)$$

By definition;
$$\pi_{ij} = \sum_{i=1}^{n} p_{ir}$$

$$\pi_{ij} = \sum_{i=1}^{n} p_{ir}$$

$$\sum_{i} \pi_{i} = \sum_{i}^{n} \sum_{r}^{n} p_{ir} = \sum_{r}^{n} \sum_{i}^{n} p_{ir} = \sum_{r}^{n} 1 = n$$
We have,
$$\sum_{i}^{N} \pi_{ii} = \sum_{r}^{N} \sum_{i}^{N} (p_{ii} is)$$

We have,
$$\sum_{j(\neq i)=1}^{n} u_{ij} = \sum_{j(\neq i)}^{n} \sum_{s(\neq i)=1}^{N} (p_{ij}J3)$$
$$= \sum_{j(\neq i)}^{N} \sum_{r=1}^{N} \sum_{s(\neq r)=1}^{N} p_{ir} p_{js/ir} [\because p(AB) = p(A) (B/A)]$$
$$= \sum_{r=1}^{N} p_{ir} \sum_{s(\neq r)=1}^{N} \sum_{j(\neq i)=1}^{N} p_{js/ir}$$
$$= \sum_{r=1}^{N} p_{ir} \sum_{s(\neq r)=1}^{N} 1$$
$$= \pi_{i}(n-1)$$
$$\sum_{i=1}^{N} \sum_{j(\neq i)=1}^{N} \pi_{ij} = n(n-1)$$
$$= \sum_{i=1}^{N} \pi_{i} (n-1)$$
$$= n (n = 1)$$
$$v(\overline{Y_{HT}}) = \sum_{i=1}^{N} \frac{Y_{i}^{2}(1-\pi_{i})\pi_{i}}{N^{2}\pi_{i}^{2}} + \sum_{j\neq i}^{N} \frac{Y_{i}Y_{j}(\pi_{ij}-\pi_{i}\pi_{j})\pi_{i}}{N^{2}\pi_{i}\pi_{j}}$$
.....(1)

This value put in (1)

$$N^{2} v(\overline{Y_{HT}}) = \sum_{i=1}^{N} \frac{Y_{i}^{2}}{\pi_{i}^{2}} \sum_{j \neq i}^{N} (\pi_{i} \pi_{j} - \pi_{ij}) - \sum_{j \neq i}^{N} \frac{Y_{i} Y_{j} (\pi_{i} \pi_{j} - \pi_{ij})}{\pi_{i} \pi_{j}}$$

$$= \sum_{j \neq i}^{N} (\pi_{i} \pi_{j} - \pi_{ij}) \left(\sum_{i=1}^{N} \frac{Y_{i}^{2}}{\pi_{i}^{2}} - \frac{Y_{i} Y_{j}}{\pi_{i} \pi_{j}} \right)$$

$$= \sum_{j \neq i}^{N} (\pi_{i} \pi_{j} - \pi_{ij}) \left(\frac{Y_{i}}{\pi_{i}} - \frac{Y_{j}}{\pi_{j}} \right)^{2}$$

$$V_{YG}(\overline{Y_{HT}}) = \frac{1}{N^{2}} \sum_{j \neq i}^{N} \left(\frac{Y_{i}}{\pi_{i}} \cdot \frac{Y_{j}}{\pi_{j}} \right)^{2} (\pi_{i} \pi_{j} - \pi_{ij})$$

$$V_{YG}(\overline{Y_{HT}}) = \frac{1}{N^{2}} \sum_{j \neq i}^{N} \delta_{i} \delta_{j} \frac{\pi_{i} \pi_{j} - \pi_{ij}}{\pi_{ij}}$$

$$E \left(V_{yG}(\overline{Y_{HT}}) \right) = \frac{1}{N^{2}} \sum_{j \neq i}^{N} E \left(\delta_{i} \delta_{j} \right) \frac{\pi_{i} \pi_{j} - \pi_{ij}}{\pi_{ij}}$$

$$= \frac{1}{N^{2}} \sum_{j \neq i}^{N} \pi_{ij} \frac{(\pi_{i} \pi_{j} - \pi_{ij})}{\pi_{ij}}$$

** The Major Draw Back of H.T Estimator are:-

- (1) Estimate of variance can assume negative value and have no inference can be
- (2) Calculation of π_i , π_{ij} , are difficult. The author have suggested that the population many be deeply stratified that a sample of 2 or 3 is sufficient to represent the strata. Yotes and Grandy claimed that the estimator proposed by them is always positive, Then, illustrated it by an exp. Of 5 popes for n = 2, Then claim is however not true for n > 2.

So it is also taken negative value but less frequently then that of H.T. ** To show that $\hat{V}_{yG}(\overline{Y_{HT}})$ is always positive for n = 2.

We know that;
$$\pi_i = p_i \left(1 + s - \frac{p_i}{1 - p_i} \right)$$
 and $\pi_{ij} = p_i p_j \left(\frac{1}{1 - p_i} + \frac{1}{1 - p_j} \right)$
Now $\pi_i = p_i \left(1 + \frac{p_i}{1 - p_i} + A \right)$ where $A = \sum_{k \neq ij}^N \frac{p_k}{1 - p_k}$ $s = A + \frac{p_i}{1 - p_i} + \frac{p_j}{1 - p_j}$

$$\begin{split} \pi_{j} &= p_{j} \bigg(1 + \frac{p_{i}}{1 - p_{i}} + A \bigg) \\ \pi_{i} \,\pi_{j} &= p_{i} p_{j} \bigg(1 + \frac{p_{j}}{1 - p_{j}} + A \bigg) \bigg(1 + \frac{p_{i}}{1 - p_{i}} + A \bigg) \\ &= p_{i} p_{j} \bigg(1 + \frac{p_{i}}{1 - p_{i}} + A + \frac{p_{j}}{1 - p_{j}} + \frac{p_{i} p_{j}}{(1 - p_{i})(1 - p_{j})} + A \frac{p_{j}}{1 - p_{j}} + A + \frac{p_{i} A}{1 - p_{i}} + A^{2} \bigg) \\ &= p_{i} p_{j} \bigg(1 + \frac{p_{i}}{1 - p_{i}} + \frac{p_{j}}{1 - p_{j}} + 2A + \frac{p_{i} p_{j}}{(1 - p_{i})(1 - p_{j})} + \bigg(\frac{p_{i}}{1 - p_{i}} + \frac{p_{j}}{1 - p_{j}} \bigg) A + A^{2} \bigg) \\ \pi_{i} \pi_{j} - \pi_{ij} &= p_{i} p_{j} \bigg(1 + \frac{p_{i}}{1 - p_{i}} + \frac{p_{j}}{1 - p_{j}} + 2A + \frac{p_{i} p_{j}}{(1 - p_{i})(1 - p_{j})} + \bigg(\frac{p_{i}}{1 - p_{i}} + \frac{p_{j}}{1 - p_{j}} \bigg) A + A^{2} \bigg) \\ &- p_{i} p_{j} \bigg(\frac{1}{1 - p_{i}} + \frac{1}{1 - p_{j}} \bigg) \end{split}$$

$$= p_i p_j \left(1 + \frac{p_i}{1 - p_i} - \frac{1}{1 - p_i} + \frac{p_j}{1 - p_j} - \frac{1}{1 - p_j} + 2A + \frac{p_i p_j}{(1 - p_i)(1 - p_j)} + \left(\frac{p_i}{1 - p_i} + \frac{p_j}{1 - p_j}\right)A + A^2 \right)$$

$$= p_i p_j \left(1 - \frac{1 - p_i}{1 - p_i} - \frac{1 - p_j}{1 - p_j} + 2A + \frac{p_i p_j}{(1 - p_i)(1 - p_j)} + \left(\frac{p_i}{1 - p_i} + \frac{p_j}{1 - p_j}\right)A + A^2 \right)$$

$$= p_i p_j \left(1 - 1 - 1 + 2A + \frac{p_i p_j}{(1 - p_i)(1 - p_j)} + \left(\frac{p_i}{1 - p_i} + \frac{p_j}{1 - p_j}\right)A + A^2 \right)$$
but, min $A = \frac{N - 2}{N - 1}$ $\left(\because \frac{(N - 2)1/N}{1 - 1/N} \right)$
Consider, $= \frac{2(N - 2)}{N - 1} - 1$

$$= \frac{2N - 4 - N + 1}{N - 1} = \frac{N - 3}{N - 2} \ge 0$$
 if N < 3.

** Midzuno Scheme of Sampling:-

Select the first unit of the sample with PPS and the remaining (n-1) units with SRSWOR.

 $\pi_i(n)$ = Probability of selecting ith unit at the 1st draw + probability of selecting at any of (n-1) draw

$$= p_i + (1 - p_i) \frac{n - 1}{N - 1}$$

$$=\frac{N-n}{N-1}p_i+\frac{\mu-1}{N-1}$$

 $\pi_{ij}(n)$ = Probability of selecting ith first draw and jth at any of the (n-1) draws + probability jth is selecting the first draw p ith included at any of (n-1) draws + neither ith nor jth unit is selected at the first draw but are included at remaining (n-1) draw

$$\begin{split} &= p_i \frac{n-1}{N-1} + p_j \frac{n-1}{N-1} + (1-p_i-p_j) \frac{n-1}{N-1} \frac{n-2}{N-2} \\ &\quad - \left(\frac{n-1}{N-1}\right) \left(p_i + p_j\right) + \left(\frac{n-1}{N-1}\right) \left(\frac{n-2}{N-2}\right) - (p_i + p_j) \left(\frac{n-1}{N-1}\right) \left(\frac{n-2}{N-2}\right) \\ &= \left(p_i + p_j \left(\frac{n-1}{N-1}\right) \left(\frac{N-n}{N-2}\right) + \left(\frac{n-1}{N-1}\right) \left(\frac{n-2}{N-2}\right) \\ &\pi_{ijk}(n) = \frac{N-n}{N-3} \frac{n-2}{N-2} \frac{n-1}{N-1} \left(p_i + p_j + p_k\right) + \left(\frac{n-1}{N-1}\right) \left(\frac{n-2}{N-2}\right) \left(\frac{n-3}{N-3}\right) \\ &\pi_{i1,i2,\dots,ir}(n) = \frac{N-n}{N-r} \frac{n-r+1}{N-r+1} \dots \frac{n-1}{N-1} \left(p_{i1} + p_{i2} + \dots + p_{ir}\right) + \frac{(n-1)\dots(n-r)}{(N-1)\dots(N-r)} \end{split}$$

Thus in this sampling scheme the probability of selection of units in the sample can always be made propertinoned to the sum of their sizes.

The another property of this sampling scheme is that Yates-Grandy is estimator or Horwity-Thompson.

Estimator's variance is always positive. They can be proved so follows,

We know that Yot's estimators of $V(\overline{Y_{HT}})$ will always be positive if $\pi_i \pi_j - \pi_{ij} > 0$, k_{ij}

$$\begin{aligned} \pi_i \pi_j - \pi_{ij} &= \left(\frac{N-n}{N-1} p_i + \frac{n-1}{N-1}\right) \left(\frac{N-n}{N-1} p_j + \frac{n-1}{N-1}\right) - \frac{N-n}{N-2}, \frac{n-1}{N-1} (p_i + p_j) \\ &- \frac{n-1}{N-1}, \frac{n-2}{N-2} \end{aligned}$$

$$= \left(\frac{N-n}{N-1}\right)^2 p_i p_j + \left(\frac{N-n}{N-1}\right) p_i \left(\frac{n-1}{N-1}\right) + \left(\frac{n-1}{N-1}\right) \left(\frac{N-n}{N-1}\right) p_j + \left(\frac{n-1}{N-1}\right) \left($$

$$= \left(\frac{N-n}{N-1}\right)^2 p_i p_j + \left(\frac{N-n}{N-1}\right) \left(\frac{n-1}{N-1}\right) (p_i + p_j) - \left(\frac{N-n}{N-2}\right) \left(\frac{n-1}{N-1}\right) (p_i + p_j) + \left(\frac{n-1}{N-1}\right) \left(\frac{n-1}{N-1} - \frac{n-2}{N-2}\right)$$

$$\begin{split} &= \left(\frac{N-n}{N-1}\right)^2 \ p_i p_j + \left(\frac{n-1}{N-1}\right) (p_i + p_j) - \left(\frac{N-n}{N-1} - \frac{N-n}{N-2}\right) + \left(\frac{n-1}{N-1}\right) \left(\frac{n-1}{N-1} - \frac{n-2}{N-2}\right) \\ &= \left(\frac{N-n}{N-1}\right)^2 \ p_i p_j + \left(\frac{n-1}{N-1}\right) (p_i + p_j) \left(\frac{-N+n}{(N-1)(N-2)}\right) + \left(\frac{n-1}{N-1}\right) \left(\frac{n-1}{N-1} - \frac{n-2}{N-2}\right) \\ &= \left(\frac{N-n}{N-1}\right)^2 \ p_i p_j - \left(\frac{n-1}{N-1}\right) \left[(p_i + p_j) \ \frac{N-n}{(N-1)(N-2)} - \frac{N-n}{(N-1)(N-2)} \right] \\ &= \frac{N-n}{(N-1)^2} \left[(N-n) \ p_i p_j - \frac{n-1}{N-2} (p_i + p_j) - \frac{n-1}{N-2} \right] > 0 \\ &\therefore \qquad (N-n) \ p_i p_j - \frac{n-1}{N-2} (p_i + p_j) - \left(\frac{n-1}{N-2}\right) > 0 \\ &\therefore \qquad \pi_i \pi_j > \pi_{ij} \end{split}$$

**** RATIO ESTIMATION ****

In sample random sampling we considered estimators using observed values of characteristic under study. Many a time the characteristic Y under study is closely related to a auxiliary characteristic X and data on X are either readily available or can be easily collected for all the unit in the population. In such situations, it is customary to consider estimators of \overline{Y} that are the data on X and are more efficient then the estimators which use data on the characteristic Y alone fact that, the data on the auxiliary variable can be used data at a later stage after selecting the sample, encourages such processors two type of these commonly used methods are

(1) The ratio type methods.

(2) The regression methods.

**** NOTATION ****

 Y_i = the value of Y, the characteristic under study for the ith unit in the population. i = 1, 2,, N

 X_i = the value of X, the auxiliary characteristic for the ith unit in the population.

 $\hat{R} = \overline{Y} / \overline{X}$: The ratio of population means.

 $\hat{R} = \overline{y} / \overline{x}$: The ratio of sample means.

The ratio estimators of the population mean \overline{Y} is that defined as $\overline{y_r} = \hat{R}\overline{X}$

Example:-Let Y = number of bullocks on a holding.

X = total area in acres

The ratio \hat{R} is an estimator of the number of bullocks for a car for holding in the population, the product of \hat{R} with \overline{X} . The coverage size of a holding in acres would prosed in estimator of \overline{Y} –1 the coverage number of bullocks per holding in the population.

y = population of Rajkot in 1992.

x = population of Rajkot in 1987.

X = total population of Rajkot in 1987.

Let us have a SRSWOR of n households for 1987 and hence get the sample total y. Then the ratio estimator of the total population of Rajkot in 1992 us given by,

$$\hat{y}_R = \frac{y}{x}X$$

**** BIAS OF RATIO ESTIMATORS:-**

The ratio estimator of the population mean \overline{Y} is given by,

$$E(e_0e_1) = E\left[\frac{\overline{y} - \overline{Y}}{\overline{Y}} \frac{\overline{x} - \overline{X}}{\overline{X}}\right] = \frac{\operatorname{cov}(\overline{y}, \overline{x}) - \overline{X}}{\overline{Y}\overline{X}} = \frac{1}{\overline{Y}\overline{X}} \frac{N - n}{Nn} S_{yx}^2$$

where $S_y^2 = \frac{1}{N - 1} \sum (y_i - \overline{Y})^2$
 $S_x^2 = \frac{1}{N - 1} \sum (x_i - \overline{X})^2$
 $S_{yx}^2 = \frac{1}{N - 1} \sum (x_i - \overline{X})(y_i - \overline{Y})$

From (ii) $y = Y(1 + e_0)$

by,

From (iii)
$$\overline{x} = \overline{X}(1+e_1)$$
 put this value in equation (i)
 $\overline{y}_R = \frac{\overline{Y}(1+e_0)}{\overline{X}(1+e_1)}\overline{X} = \overline{y}\left(\frac{1+e_0}{1+e_1}\right) = \overline{y}(1+e_0)(1+e_1)^{-1}$
 $= \overline{y}(1+e_0)(1-e_1+e_1^2)$ [In the expansion s terms having than power 3 are negleted
 $= \overline{y}(1+e_0-e_1-e_0e_1+e_1^2)$

Subtract \overline{Y} from both the side in above expression we get,

$$\overline{Y}_R - \overline{Y} = \overline{Y} (e_0 - e_1 - e_0 e_1 + e_1^2)$$

Bias of the ratio estimator to the first degree of approximation is given

Bias
$$(\overline{y_R}) = E[(\overline{y_R}) - \overline{y}] = E(\overline{Y_R}) - \overline{Y}$$

$$= \overline{Y} E(e_0 - e_1 - e_0 e_1 + e_0^2)$$

$$= \overline{Y} E(e_1^2 - e_0 e_1) \quad [\because E(e_0) = E(e_1) = 0]$$

$$= \overline{Y} \Big[E(e_1^2) - E(e_0 e_1) \Big]$$

$$= \overline{Y} \Big[\frac{1}{\overline{X}^2} \frac{N - n}{Nn} S_x^2 - \frac{1}{\overline{X}\overline{Y}} \frac{N - n}{Nn} S_{yx}^2 \Big]$$

$$= \overline{Y} \Big[\frac{N - n}{Nn} \left(\frac{S_x^2}{\overline{X}^2} \frac{S_{yx}^2}{\overline{Y}\overline{X}} \right) \Big]$$

$$= \overline{Y} \frac{1 - f}{n} (CX^2 - \tau C_y C_x) \qquad \left[\because \frac{n}{N} = f, \ S_{yx} = \tau S_y S_x, \ \tau = S_{yx} / S_y S_x \right]$$

where, $\tau = correlation$ coefficient between y and x.

 $C_x = S_x/\overline{X}$ = coefficient of variation of x. $C_y = S_y/\overline{Y}$ = coefficient of variation of y. f = n/N sample fraction.

REMARK:-

Ratio estimator is unbiased if

bias $(\overline{Y_R}) = 0$. i.e. if $\overline{Y} \frac{1-f}{n} \left[\frac{S_x^2}{\overline{X}^2} - \tau \frac{S_y S_x}{\overline{YX}} \right] = 0$ i.e. if $\frac{S_x^2}{\overline{X}^2} - \frac{S_y S_x}{\overline{YX}} = 0$ i.e. if $\frac{S_x^2}{\overline{X}^2} = \frac{S_{yx}}{\overline{YX}}$ i.e. if $\overline{Y} = \frac{S_{yx}}{S_x^2} \overline{X}$ i.e. if $\overline{Y} = \beta \overline{X}$, $\beta = \frac{S_{yx}}{S_x^2} = regression coefficient of y on x$

i.e. if the line of regression of y on x is a straight line passing through origin.



**** MEAN SQUARE ERROR OF RATIO ESTIMATOR :-**

The mean square error of ratio estimator of the populoation mean is given by,

$$MSE(\overline{Y_{R}}) = E\left[\overline{Y_{R}} - \overline{Y}\right]^{2}$$

$$= E\left[\overline{Y}(e_{0} - e_{1})\right]^{2}$$

$$= \overline{Y}^{2}E(e_{0}^{2} - 2e_{0}e_{1} + e_{1}^{2})$$

$$= \overline{Y}^{2}\left[E(e_{0}^{2}) - 2E(e_{0}e_{1}) + E(e_{1}^{2})\right]$$

$$= \overline{Y}^{2}\frac{N - n}{Nn}\left[\frac{S_{y}^{2}}{\overline{Y}^{2}} - \frac{2S_{yx}}{\overline{Y}\overline{X}} + \frac{S_{x}^{2}}{\overline{X}^{2}}\right]$$

$$= \frac{N - n}{Nn}\left[S_{y}^{2} - 2S_{yx}\frac{\overline{Y}}{\overline{X}} + S_{x}^{2}\frac{\overline{Y}^{2}}{\overline{X}^{2}}\right]$$

$$= \frac{N-n}{Nn} \left[S_{y}^{2} - 2RS_{yx} + R^{2}S_{x}^{2} \right] \qquad \left[\because R = \frac{\overline{Y}}{\overline{X}} \right]$$

OR
$$MSE(\overline{y_{R}}) = \overline{Y}^{2} \ \frac{1-f}{n} \left[C_{y}^{2} - 2\rho C_{y}C_{x} + C_{x}^{2} \right]$$

**** COMPARISON OF RATIO ESTIMATOR WITH SRSWOR:-**

In SRSWOR, we have $V(\overline{y_n})_{SRS} = \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 = \frac{1-f}{n} S_y^2$ (:: f = n/N) For ratio estimator we have, $MSE\left(\overline{y_{R}}\right) = \frac{1-f}{n} \left[S_{y}^{2} - 2RS_{yx} + R^{2}S_{x}^{2}\right]$ Ratio estimator is better then SRSWOR if $V(\overline{y_n})_{SRS} - MSE(\overline{y_R}) > 0$ $\frac{1-f}{n} \left[S_{y}^{2} \right] - \frac{1-f}{n} \left[S_{y}^{2} - 2RS_{yx} + R^{2}S_{x}^{2} \right] > 0$ i.e. if $\frac{1-f}{n} \left[2RS_{yx} + R^2 S_x^2 \right] > 0$ i.e. if $2RS_{yx} > R^2 S_x^2$ $2R\tau S_{y}S_{x} > R^{2}S_{x}^{2}$ i.e. if $\rho > \frac{1}{2} \frac{S_x}{S_y} R > \frac{1}{2} \frac{S_x}{S_y} \frac{\overline{Y}}{\overline{X}}$ $(:: -1 < \rho < 1)$ $\rho > \frac{1}{2} \frac{C_x}{C_y}$ i.e. if $\rho > \frac{1}{2}$ when $C_x = C_y$ i.e. if

**** ESTIMATION OF VARIANCE OF RATIO ESTIMATOR:-**

The estimate of the variance of ratio estimator is given by,

$$EST(\overline{y_{R}}) = \frac{1-f}{n} \left[s_{y}^{2} - 2rs_{yx} + r^{2}s_{x}^{2} \right]$$

where $r = \overline{y}/\overline{x}$
 $s_{y}^{2} = 1/n - 1 \sum_{i=1}^{n} (y_{i} - \overline{y})^{2}$
 $s_{x}^{2} = 1/n - 1 \sum_{i=1}^{n} (x_{i} - \overline{x})^{2}$

$$s_{yx}^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (y_{i} - \overline{y}) (x_{i} - \overline{x})$$
Now,

$$Est \ V(\overline{y}_{R}) = \frac{1-f}{n} \frac{1}{n-1} \left[\sum_{i=1}^{n} (y_{i} - \overline{y})^{2} - 2r \sum_{i=1}^{n} (y_{i} - \overline{y}) (x_{i} - \overline{x}) + r^{2} \sum_{i=1}^{n} (x_{i} - \overline{x})^{2} \right]$$

$$= \frac{1-f}{n} \frac{1}{n-1} \left[\sum_{i=1}^{n} \{ (y_{i} - \overline{y}) - r(x_{i} - \overline{x}) \}^{2} \right]$$

$$= \frac{1-f}{n} \frac{1}{n-1} \left[\sum_{i=1}^{n} \{ y_{i} - rx_{i} \}^{2} \right]$$

$$= \frac{1-f}{n} \frac{1}{n-1} \left[\sum_{i=1}^{n} y_{i}^{2} - 2r \sum_{i=1}^{n} y_{i}x_{i} + r^{2} \sum_{i=1}^{n} x_{i}^{2} \right]$$

**** UNBIASED RATIO TYPE ESTIMATOR:-**

An Unbiased Ratio Type Estimator of the population mean \overline{Y} is defined

as,

Now, taking expectation on both sides,

$$E[\overline{Y}_{GH}] = E[\overline{r} \ \overline{x}] + \frac{n(N-1)}{N(n-1)} E(\overline{y} - \overline{r} \ \overline{x})$$

$$= \overline{R} \ \overline{X} + \frac{N-1}{N} E\left[\frac{n\overline{y} - n\overline{r} \ \overline{x}}{n-1}\right]$$

$$= \overline{R} \ \overline{X} + \frac{N-1}{N} E\left[\frac{\sum_{i=1}^{n} y_i - n\overline{r} \ \overline{x}}{n-1}\right] \qquad [\because \sum y_i/n = \overline{Y}]$$

$$= \overline{R} \ \overline{X} + \frac{N-1}{N} E\left[\frac{\sum_{i=1}^{n} r_i x_i - n\overline{r} \ \overline{x}}{n-1}\right] \qquad [\because \sum y_i/n = \overline{Y}]$$

$$= \overline{R} \ \overline{X} + \frac{N-1}{N} \ E[s_{rx}]$$

$$= \overline{R} \ \overline{X} + \frac{N-1}{N} \ S_{rx} \qquad [\because (Es_{rx}) \cong S_{rx}]$$

$$= \overline{R} \ \overline{X} + \frac{N-1}{N} \left[\frac{\sum_{i=1}^{N} r_i x_i - N \ \overline{R} \ \overline{X}}{N-1} \right]$$

$$= \overline{R} \ \overline{X} + \frac{1}{N} \ \sum_{i=1}^{N} r_i x_i - \overline{R} \ \overline{X}$$

$$= \frac{1}{N} \ \sum_{i=1}^{N} r_i x_i$$

$$= \frac{1}{N} \ \sum_{i=1}^{N} y_i \qquad [\because r_i = y_i / x_i]$$

$$= \overline{Y}$$

$$E(\overline{y}_{GH}) = \overline{Y}$$

EXAMPLE :- Complete the variance of ratio estimators of the population mean : When N = 100, n = 25, $\rho = \frac{1}{2}$ C_x = C_y complete it with SRSWOR :-

When N = 100, n = 25,
$$\rho = \frac{1}{2}$$
 C_x = C_y complete it with SRSW

$$MSE(\overline{y_R}) = \overline{Y}^2 \frac{1-f}{n} [C_y^2 - 2\rho C_y C_x + C_x^2]$$

$$= \overline{Y}^2 \frac{1-f}{n} [C_y^2 - C_y^2 2\rho + C_y^2]$$

$$= \overline{Y}^2 \frac{1-f}{n} 2C_y^2 [1-\rho]$$

$$= \overline{Y}^2 \frac{1-f}{n} \frac{2S_y^2}{\overline{Y}^2} [1-\rho] \qquad [\because C_y^2 = S_y^2/\overline{Y}^2]$$

$$f = \frac{n}{N} = \frac{25}{100} = \frac{1}{4} = 0.25$$

$$= \frac{0.75}{25} \times 2. S_y^2 (1-\frac{1}{2})$$

$$= \frac{75}{25} \frac{S_y^2}{100}$$

$$MSE(\overline{y_R}) = 0.03 S_y^2$$

For SRSWOR

$$V(\bar{y}_n) = \frac{N-n}{Nn} S_y^2$$

= $\frac{100-25}{2500} S_y^2$
= $0.03 S_y^2$

**** REGRESSION METHOD OF ESTIMATOR:-**

The difference estimator of the population mean \overline{Y} is defined as, $\overline{Y}_d = \overline{y} + d(\overline{X} - \overline{x})$ (*) where d is some constant Now, $E(\overline{Y}_d) = E(\overline{y}) + d E(\overline{X} - \overline{x})$

$$= \overline{y} + d(\overline{X} - \overline{X})$$
$$= \overline{y}$$

i.e. \overline{y}_d is an unbiased estimator of \overline{y} next,

$$V(\overline{y}_d) = E[\overline{y}_d - E(\overline{y}_d)]^2$$

= $E[\overline{y}_d - \overline{y}]^2$
= $E[\overline{y} + d(\overline{X} - \overline{x}) - \overline{y}]^2$ [from (*)]
= $E[(\overline{y} - \overline{Y}) + d(\overline{X} - \overline{x})]^2$
= $E[(\overline{y} - \overline{Y}) - d(\overline{x} - \overline{X})]^2$
= $E[(\overline{y} - \overline{Y})^2 - 2dE(\overline{y} - \overline{Y})(\overline{x} - \overline{X}) + d^2E(\overline{x} - \overline{X})^2$
= $V(\overline{y}) - 2d \operatorname{cov}(\overline{y}, \overline{x}) + d^2V(\overline{x})$

To determine constant d we minimize $V(\bar{y}_d)$ w.r.t. d.

$$\frac{\partial V(\bar{Y}_d)}{\partial d} = 0 \Rightarrow \frac{\partial}{\partial d} \Big[V(\bar{y}) - 2d \operatorname{cov}(\bar{y}, \bar{x}) + d^2 V(\bar{x}) \Big] = 0$$

$$\Rightarrow -2 \operatorname{cov}(\bar{y}, \bar{x}) + 2d V(\bar{x}) = 0$$

$$\Rightarrow 2 \operatorname{cov}(\bar{y}, \bar{x}) = 2d V(\bar{x})$$

$$\Rightarrow \operatorname{cov}(\bar{y}, \bar{x}) = d V(\bar{x})$$

$$\Rightarrow d = \frac{\operatorname{cov}(\bar{y}, \bar{x})}{V(\bar{x})}$$

$$= \frac{1 - f}{n} \frac{S_{yx}}{S_x^2} = \frac{S_{yx}}{S_x^2} = \beta$$

 β is Regression coefficient of y over x.

 \therefore the difference estimators becomes

$$\bar{y}_{\beta} = \bar{y} + \beta(\bar{X} - \bar{x})$$

Where β the population regression coefficient of y on x is unknown it will be replaced by simple regression coefficient of y on x, b-we get the linear regression estimator of \overline{Y} .

i.e. $\overline{Y}_{l_r} = \overline{y} + b(\overline{X} - \overline{x})$ where $b = s_{yx}/s_x^2$, lr = linear regression.

**** BIAS OF REGRESSION ESTIMATOR:-**

Let us define, (1)
$$e_0 = \frac{\overline{y} - \overline{Y}}{\overline{Y}}$$
 $E(e_0) = 0$
(2) $e_1 = \frac{\overline{x} - \overline{X}}{\overline{X}}$ $E(e_1) = 0$
(3) $e_2 = \frac{b - \beta}{\beta}$ $E(e_2) = 0$
 $E(e_0^2) = \frac{V(\overline{y})}{\overline{Y}^2} = \frac{1}{\overline{Y}^2} \frac{N - n}{Nn}$ $S_y^2 = \frac{1 - f}{n} S_y^2 / \overline{Y}^2$
 $E(e_1^2) = \frac{1 - f}{n} S_x^2 / \overline{X}^2$
 $E(e_0e_1) = \frac{1 - f}{n} E\left[\frac{(\overline{y} - \overline{Y})(\overline{x} - \overline{X})}{\overline{Y} \, \overline{X}}\right] = \frac{\operatorname{cov}(\overline{y}, \overline{x})}{\overline{Y} \, \overline{X}} = \frac{1 - f}{n} \frac{S_{yx}}{\overline{x} \, \overline{y}}$
 $E(e_1e_2) = \frac{\operatorname{cov}(\overline{x}, b)}{\overline{X} \, \beta}$

There fore the regression estimator becomes.

$$\begin{split} \overline{Y}_{lr} &= \overline{y} + b \left(\overline{X} - \overline{x} \right) \\ &= \overline{y} (1 + e_0) + (e_2 + 1) \beta \left(-e_1 \overline{X} \right) \\ &= \overline{y} (1 + e_0) - \beta \overline{X} e_1 (1 + e_2) \end{split}$$

i.e. $\overline{Y}_{lr} - \overline{Y} = \overline{y}e_0 - \beta \overline{X} e_1(1+e_2)$ The bias of regression estimator is

The bias of regression estimator is given by:

$$Bias (\overline{Y}_{lr}) = E[\overline{Y}_{lr} - \overline{Y}]$$

$$= E[\overline{Y}_{e_0} - \beta \overline{X} \ e_1(1 + e_2)]$$

$$= E[\overline{Y}_{e_0} - \beta \overline{X} \ (e_1 + e_1e_2)]$$

$$= \overline{Y} \ E(e_0) - \beta \overline{X} \ E(e_1 + e_1e_2)$$

$$= \overline{Y} \ 0 - \beta \overline{X} \ E(e_1) - \beta \overline{X} \ E(e_1e_2)$$

$$= 0 - 0 - \beta \overline{X} \ E(e_1e_2)$$

$$= -\beta \overline{X} \ E(e_1e_2)$$

$$= -\beta \overline{X} \ \frac{\text{cov} \ (\overline{x}, b)}{\beta \ \overline{X}}$$

Bias $(\overline{Y}_{lr}) = -\cot(\overline{x}, b)$

**** MEAN SQUARE ERROR OF REGRESSION ESTIMATOR:-**

Mean square error of the regression estimator of the population mean to the first degree of approximation is given by:

$$\begin{split} MSE(\bar{Y}_{lr}) &= E\left[\overline{Y}_{lr} - \overline{Y}\right]^{2} \\ &= E\left[\overline{Y}e_{0} - \beta\overline{X} \ e_{1}(1 + e_{2})\right]^{2} \\ &= E\left[\overline{Y}^{2}e_{0}^{2} - 2\beta\overline{X}\overline{Y} \ e_{1}e_{0}(1 + e_{2}) + \beta^{2}\overline{X}^{2} \ e_{1}^{2}(1 + e_{2})^{2}\right] \\ &= \overline{Y}^{2}E(e_{0}^{2}) - 2\beta\overline{X}\overline{Y}E(e_{1}e_{0})\frac{b}{\beta} + \beta^{2}\overline{X}^{2} \ e_{1}^{2}\frac{b^{2}}{\beta^{2}} \qquad \left[\because 1 + e_{2} = \frac{b}{\beta}\right] \\ &= \overline{Y}^{2}E(e_{0}^{2}) - 2\overline{X}\overline{Y}bE(e_{1}e_{0}) + b^{2}\overline{X}^{2} \ E(e_{1}^{2}) \\ &= \overline{Y}^{2}\frac{1 - f}{n}\frac{S_{y}^{2}}{\overline{Y}^{2}} - 2\overline{X}\overline{Y}b\frac{1 - f}{n}\frac{S_{yx}}{\overline{Y}\overline{X}} + b^{2}\overline{X}^{2}\frac{1 - f}{n}\frac{S_{x}^{2}}{\overline{X}} \\ &= \frac{1 - f}{n}\left[S_{y}^{2} - 2bS_{yx} + b^{2}S_{x}^{2}\right] \qquad \left[\because S_{yx} = \rho S_{y}S_{x}\right] \\ &= \frac{1 - f}{n}\left[S_{y}^{2} - 2\beta S_{y}S_{x} + b^{2}S_{x}^{2}\right] \qquad \left[b = \frac{S_{yx}}{S_{x}^{2}} = \frac{\rho S_{y}S_{x}}{S_{x}^{2}}\right] \\ &= \frac{1 - f}{n}\left[S_{y}^{2} - 2\rho^{2}S_{y}^{2} + \rho^{2}S_{y}^{2}\right] \\ &= \frac{1 - f}{n}\left[S_{y}^{2} - 2\rho^{2}S_{y}^{2} + \rho^{2}S_{y}^{2}\right] \end{split}$$

Which is mean square error of regression estimator.

** COMPARISON OF REGRESSION ESTIMATOR WITH SRSWOR AND RATIO ESTIMATOR:-

In SRSWOR we have $V(\overline{Y}_n) = \frac{1-f}{n}S_y^2$ For regression estimation we have $MSE(\hat{Y}_{lr}) = \frac{1-f}{n}S_y^2(1-\rho^2)$ now, $V(\overline{Y}_n) - MSE(\hat{Y}_{lr}) = \frac{1-f}{n}S_y^2 - \frac{1-f}{n}S_y^2(1-\rho^2)$ $= \frac{1-f}{n}S_y^2 \rho^2 > 0$

i.e.
$$V(\overline{Y}_n) > MSE(\hat{Y}_{lr})$$

i.e. Regression estimator is superior to SRSWOR for ratio estimator, we have,

$$MSE \ (\hat{\bar{Y}}_{R}) = \frac{1-f}{n} \left[S_{y}^{2} - 2R\rho S_{y} S_{x} + R^{2} S_{x}^{2} \right]$$

now

$$MSE \ (\hat{Y}_{R}) - MSE \ (\overline{Y}_{lr}) = \frac{1 - f}{n} \Big[S_{y}^{2} - 2R\rho S_{y} S_{x} + R^{2} S_{x}^{2} \Big] - \frac{1 - f}{n} \Big[S_{y}^{2} (1 - \rho)^{2} \Big] \\= \frac{1 - f}{n} \Big[S_{y}^{2} - 2R\rho S_{y} S_{x} + R^{2} S_{x}^{2} - S_{y}^{2} + \rho^{2} S_{y}^{2} \Big] \\= \frac{1 - f}{n} \Big[R^{2} S_{x}^{2} - 2R\rho S_{y} S_{x} + \rho^{2} S_{y}^{2} \Big] \\= \frac{1 - f}{n} \Big[RS_{x} - \rho S_{y} \Big]^{2} \ge 0$$

i.e. $MSE(\hat{\overline{Y}_{R}}) \ge MSE(\overline{Y}_{lr})$

i.e. regression estimation is superior to ratio estimation and equally efficient.

If
$$RS_x - \rho S_y = 0$$

i.e. If $RS_x = \rho S_y$
i.e. if $\frac{\overline{Y}}{\overline{X}} S_x = \rho S_y$
i.e. if $\overline{Y} = \rho \frac{S_y}{S_x} \overline{X}$ $\left| \because \beta = \frac{S_{yx}}{S_x^2} = \frac{\rho S_y S_x}{S_x^2} = \frac{S_y}{S_x} \right|$
i.e. if $\overline{Y} = \beta \overline{X}$
i.e. if $\overline{Y} = \beta \overline{X}$

i.e. if the line of regression of Y on X is a straight line passing through origin.

<u>NOTE:</u> When \overline{Y}_{lr} and $\hat{\overline{Y}}_{R}$ are equally efficient there use $\hat{\overline{Y}}_{R}$ because it dose not require the compitation of \mathbf{b}_{i} .

**** APPLICATION:-**

When the measurement of one variable (study variable) is complicative costly and time consuming where as the measurement of another variable (auxiliary variable) is simple , cheaper and quicker then the regression estimator is used.

Suppose, We wish to estimate the average area of a leaf on a certain plants the weight of a leaf. Here the measurement of y is complicative, costly and time consuming where as measurement of x is simple, cheaper and quicker.

We can weight all the leaves together and get the total $\sum X = N\overline{X}$. Let us have a SRSWOR of n leaves and hence the value of y and x we compute $\overline{Y}, \overline{X}$ and b for the sample then the regression estimator of the average of a leaf is given by,

$$\overline{Y}_{lr} = \overline{y} + b(\overline{X} - \overline{x})$$

**** SYSTEMETIC SAMPLING ****

In systematic sampling the 1st sample is selected according to given pattern. Suppose, there are N units in population. Let n be a sample size. Assume that N = nk, We list the units of population. We draw a random number term 1 to k. Suppose it is i (i = 1, 2, ..., k.) then the ith unit is selected in the sample the remaining (n-1) units are selected as the units having no.

 $i+k, i+2k, i+3k, \dots, i+(n-1)k.$

**** SYSTEMETIC SAMPLING ****

In systematic sampling the 1st sample is selected according to given pattern. Suppose, there are N units in population. Let n be a sample size. Assume that N = nk, We list the units of population. We draw a random number term 1 to k. Suppose it is i (i = 1, 2, ..., k.) then the ith unit is selected in the sample the remaining (n-1) units are selected as the units having no.

 $i+k, i+2k, i+3k, \dots, i+(n-1)k.$

thus the systematic sample of size n the units with no consist of i, i+k, i+2k, i+3k,.....i+(n-1)k. Where k is called the interval of systematic sample. This is called the linear systematic sampling.

e.g. N = 100, n = 10, k = N/n = 10.

The 1st unit is selected by drawing a random no. from 1 to 10. Suppose it is 7. Then 7th unit is selected in the sample. Thus, systematic sample of size 10 contains the units as 7,17, 27, 37, 47, 57, 67, 77, 87, 97.

**** SELECTION OF ALL SYSTEMATIC SAMPLE:-**

 $1 \qquad 2 \qquad \cdots \qquad i \qquad \cdots \qquad k$ $1+k \qquad 2+k \qquad \cdots \qquad i+k \qquad \cdots \qquad 2k$ $1+2k \qquad 2+2k \qquad \cdots \qquad i+2k \qquad \cdots \qquad 3k$ $\vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots$ $1+(n-1)k \qquad 2+(n-1)k \qquad \cdots \qquad i+(n-1)k \qquad \cdots \qquad (n-1)k$

1) There k systematic sample each of size n

2) Each systematic sample has probability 1/k of being selected.

**** RELATION WITH STRATIFIED SAMPLING:-**

Systematic sampling is similar to Stratified sampling in which N units are grouped into n strata each of size k and one unit per stratum is selected. But the difference between them lies in fact that the unit occupy the same relative position in different strata where as in stratified sampling units are selected at random in different strata.

*** VALUE OF VARIABLE Y IN THE POPULATION:-**

1^{st}	2^{nd}	3 ^{<i>rd</i>}	•••	i^{th}	•••	n^{th}	mean
Y_{11}	<i>Y</i> ₂₁	Y_{31}	•••	Y_{i1}	•••	Y_{n1}	\overline{Y} 1
<i>Y</i> ₁₂	Y_{22}	<i>Y</i> ₃₂	•••	Y_{i2}	•••	Y_{n2}	\overline{Y} 2
:	÷	÷	÷	:	÷	÷	:
Y_{1j}	Y_{2j}	Y_{3j}		Y_{ij}		Y_{nj}	\overline{Y} j
:	:	:	÷	:	÷	:	:
Y_{1n}	Y_{2n}	Y_{3n}	•••	Y_{in}	•••	Y_{nn}	$\overline{Y} n$
mean $\overline{Y}1$	\overline{Y} 2	\overline{Y} 3	•••	$\overline{Y}i$	•••	$\overline{Y} n$	\overline{Y}

* NOTATIONS:-

Let, Y_{ij} = value of jth unit in the ith systematic sample. Where i =1,2,...k. and k = 1,2,...n. $\overline{Y}_i = \sum_{j=1}^n y_{ij}/k$ = mean of the ith systematic sample. $\overline{Y}_j = \sum_{i=1}^k y_{ij}/k$ = mean of the jth unit.

$$\overline{Y} = \frac{\sum_{i=1}^{k} \sum_{j=1}^{n} Y_{ij}}{nk} = \frac{1}{k} \sum_{j=1}^{n} \overline{Y}_{i}$$

= mean of the population.
= mean of the systematic sample mean.

**** EXPECTED VALUE OF SYSTEMATIC SAMPLE MEAN:-**

$$E(\overline{Y}_{sys}) = \sum_{i=1}^{k} \frac{1}{k} \overline{Y}_{i}$$

$$= \frac{1}{k} \sum_{i} \overline{Y}_{i}$$

$$= \frac{1}{nk} \sum_{i} \sum_{j} Y_{ij}$$

$$= \frac{1}{N} \sum_{i} \sum_{j} Y_{ij}$$

$$= \overline{Y} = \text{Population mean}$$

$$E(\overline{Y}_{sys}) = \overline{Y} = \overline{Y}_{N}$$

Systematic sample mean (\overline{Y}_{sys}) is an unbiased estimator of population mean \overline{Y}_{N} .

* Variance of Systematic Mean:-

$$V(\overline{Y}_{sys}) = E\left[\overline{Y}_{sys} - E(\overline{Y}_{sys})\right]^{2}$$
$$= E\left[\overline{Y}_{sys} - \overline{Y}..\right]^{2}$$
$$= E\left[\frac{\sum_{i}\sum_{j}Y_{ij}}{nk} - \overline{Y}..\right]^{2}$$
$$= E\left[\sum_{i}\frac{1}{k}\overline{Y}_{i} - \overline{Y}..\right]^{2}$$
$$= \sum_{i}\frac{1}{k} (\overline{Y}_{i} - \overline{Y}..)^{2}$$

(I) ALTERNATIVE EXPRESSION OF VARIANCE :-

$$V(\bar{Y}_{sys}) = \frac{N-1}{N} S_{y}^{2} - \frac{N-k}{N} S_{wsys}^{2}$$

Where $S_{y}^{2} = \frac{1}{N-1} \sum_{i} \sum_{j} (Y_{ij} - \bar{Y}_{..})^{2}$ (1)

= mean square among the units in the whole population.

=mean square among the units within the systematic sample. **Proof:-** we have,

$$\sum_{i} \sum_{j} (Y_{ij} - \overline{Y}_{..})^{2} = \sum_{i} \sum_{j} (Y_{ij} - \overline{Y}_{i} + \overline{Y}_{i} - \overline{Y}_{..})^{2}$$

= $\sum_{i} \sum_{j} (Y_{ij} - \overline{Y}_{i})^{2} + \sum_{i} \sum_{j} (\overline{Y}_{i} - \overline{Y}_{..})^{2} + PT$
i.e. $(N-1)S_{y}^{2} = k(n-1)S_{wsys}^{2} + n\sum_{i} (Y_{i} - \overline{Y}_{..})^{2}$ from (1) and (2)

i.e.
$$(N-1)S_y^2 = (N-k)S_{wsys}^2 + nk \ V(\overline{Y}_{sys})$$

i.e.
$$V(\overline{Y}_{sys}) = \frac{N-1}{N}S_y^2 - \frac{N-k}{N}S_{wsys}^2$$

*** COMPARISION WITH SRSWOR:-**

In SRSWOR we have,

$$V(\overline{Y}_n) = \frac{N-n}{Nn} S_y^2$$

In systematic sampling we have,

$$V(\overline{Y}_{sys}) = \frac{N-1}{N} S_y^2 - \frac{N-K}{N} S_{wsys}^2$$

Systematic sample is better then SRSWOR i.e If $V(\overline{Y}_n) > V(\overline{Y}_{sys})$ i.e If $\frac{N-n}{Nn} S_y^2 > \frac{N-1}{N} S_y^2 - \frac{N-k}{N} S_{wsys}^2$ i.e. If $\left(\frac{N-1}{N} - \frac{N-n}{Nn}\right) S_y^2 < \frac{N-k}{N} S_{wsys}^2$ i.e. If $\left(\frac{n-1}{n}\right) S_y^2 < \frac{k(n-1)}{nk} S_{wsys}^2$ i.e. If $S_y^2 < S_{wsys}^2$

i.e. If mean square among the units in the whole population < mean square among the units within the systematic sample.

$$(\mathbf{II}) \ V(\overline{Y}_{sys}) = \frac{N-1}{nN} \ S_y^2 \ [1+(n-1)\rho]$$

where, $\rho = \frac{\sum_{j\neq i}^n \sum_{i=1}^k (y_{ij} - \overline{Y})(y_{ij} - \overline{Y}..) / kn(n-1)}{\sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \overline{Y}..)^2 / nk}$
$$(1) \leftarrow = \frac{\sum_{i=1}^k \sum_{j\neq i=1}^n (y_{ij} - \overline{Y}..)(y_{ij} - \overline{Y}..)}{(n-1) \ (N-1) \ S_y^2}$$

= intraclass correlation coefficient between pair of units within the systematic sample.

Proof:- we have,

$$V(\overline{Y}_{sys}) = \frac{1}{k} \sum_{i=1}^{k} (\overline{Y}_{i} - \overline{Y}_{..})^{2}$$

$$n^{2}V(\overline{Y}_{sys}) = \frac{n^{2}}{k} \sum_{i=1}^{k} (\overline{Y}_{i} - \overline{Y}_{..})^{2} \qquad \text{[multiplying by n}^{2}$$

$$n^{2}kV(\overline{Y}_{sys}) = \sum_{i} (n\overline{Y}_{i} - n\overline{Y}_{..})^{2}$$

$$n^{2}kV(\overline{Y}_{sys}) = \sum_{i} \left(\sum_{j} Y_{ij} - n\overline{Y}_{..} \right)^{2} \qquad \left| \frac{\sum_{j} Y_{ij}}{n} = \overline{Y}_{i} \right|$$

$$= \sum_{i=1}^{k} \sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{..})^{2} + \sum_{i} \sum_{j \neq j=1}^{m} (Y_{ij} - \overline{Y}_{..})(Y_{ij} - \overline{Y}_{..})$$

$$= (N-1) S_{y}^{2} + \rho (n-1) (N-1) S_{y}^{2} \qquad |\because from (1) \& equation of S_{y}^{2}$$

$$n^{2}kV(\overline{Y}_{sys}) = (N-1) S_{y}^{2} [1 + (n-1) \rho]$$
i.e. $V(\overline{Y}_{sys}) = \frac{N-1}{nN} S_{y}^{2} [1 + (n-1) \rho]$

* Comparison With SRSWOR:-

In SRSWOR,

$$V(\overline{Y}_n) = \frac{N-n}{Nn} S_y^2$$
$$V(\overline{Y}_{sys}) = \frac{N-1}{nN} S_y^2 [1 + (n-1)\rho]$$

Systematic sampling is better than SRSWOR

i.e If
$$V(\overline{Y}_{sys}) < V(\overline{Y}_n)$$

i.e. if
$$\frac{N-1}{nN}S_{y}^{2}[1+(n-1)\rho] < \frac{N-n}{Nn}S_{y}^{2}$$

i.e. if $\left(\frac{N-1}{nN}-\frac{N-n}{Nn}\right)S_{y}^{2} < -\frac{N-1}{nN}S_{y}^{2}(n-1)\rho$
i.e. if $\frac{N-1}{nN}S_{y}^{2} < -\frac{N-1}{nN}S_{y}^{2}(n-1)\rho$
i.e. if $1 < -(N-1)\rho$
i.e. if $1 > (N-1)\rho$
i.e. if $\rho < \frac{1}{N-1}$

(III)
$$V(\overline{Y}_{sys}) = \frac{S_{wst}^{2}}{n} \left[\frac{N-n}{N} + (n-1) \rho_{wst} \right]$$

(i) where $S_{wst}^{2} = \sum_{i=1}^{k} \sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{j})^{2} / n(k-1)$
$$\rho_{wst} = \frac{\sum_{i=1}^{k} \sum_{j\neq i=1}^{n} (Y_{ij} - \overline{Y}_{.j})(Y_{ij} - \overline{Y}_{.j}) / kn(n-1)}{\sum_{i} \sum_{j} (Y_{ij} - \overline{Y}_{.j})^{2} / n(k-1)}$$
$$= \frac{\sum_{i=1}^{k} \sum_{j\neq i=1}^{n} (Y_{ij} - \hat{Y}_{ij})(Y_{ij} - \overline{Y}_{.j}) / kn(n-1)}{S_{wst}^{2}}$$

= Correlation coefficient between pairs of units within

the strata.

$$\begin{aligned} \mathbf{Proof:} \quad V\left(\overline{Y}_{sys}\right) &= \frac{1}{k} \sum_{i=1}^{k} (\overline{Y}_{i} - \overline{Y}_{..})^{2} \\ n^{2} k V(\overline{Y}_{sys}) &= \sum_{i=1}^{k} (n\overline{Y}_{i} - \overline{Y}_{..})^{2} \qquad |\because multiplying by n^{2} \\ &= \sum_{i=1}^{k} \left(\sum_{j=1}^{n} Y_{ij} - n\overline{Y}_{..} \right)^{2} \qquad \left| \overline{Y}_{i} &= \sum_{j} Y_{ij} / n \\ &= \sum_{i=1}^{k} \left(\sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{..}) \right)^{2} \\ &= \sum_{i=1}^{k} \left(\sum_{j=1}^{n} Y_{ij} - \sum_{i} y_{.j} \frac{n}{n} \right)^{2} \qquad \left| \overline{Y} &= \frac{\sum_{j=1}^{j} y_{.j}}{n} \right| \end{aligned}$$

$$= \sum_{i=1}^{k} \left(\sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{\cdot j}) \right)^{2}$$

$$= \sum_{i} \sum_{j} (y_{ij} - \overline{y}_{\cdot j}) + \sum_{i=1}^{k} \sum_{j \neq j'=1}^{n} (y_{ij} - y_{\cdot j})(y_{ij'} - y_{\cdot j'})$$

$$= n(k-1) S_{wst}^{2} + kn(n-1) \rho_{wst} S_{wst}^{2} \qquad |\because from(1) \& (2)$$

$$= S_{wst}^{2} [n(k-1) + kn(n-1) \rho_{wst}]$$

$$V(\overline{Y}_{sys}) = \frac{S_{wst}^{2}}{n^{2}k} [n(k-1) + kn(n-1) \rho_{wst}]$$

$$= \frac{S_{wst}^{2}}{n^{2}k} \left[\frac{n(k-1)}{nk} + \frac{kn(n-1) \rho_{wst}}{nk} \right]$$

$$= \frac{S_{wst}^{2}}{n} \left[\frac{nk-n}{nk} + (n-1) \rho_{wst} \right]$$

* Comparison With Stratified Sampling

There are N units divided into n strata each of size k and 1 unit per strata is selected.

$$\begin{split} V(\overline{Y}_{st}) &= \sum_{j=1}^{n} p_{j}^{2} \left(\frac{1}{n_{j}} - \frac{1}{N_{j}} \right) S_{j}^{2} \\ \text{Where } p_{j} &= \frac{N_{j}}{N} = \frac{k}{nk} = \frac{1}{n} = \frac{n_{j}}{n} \qquad \therefore n_{j} = 1, \quad N_{j} = k \\ S_{j}^{2} &= \frac{1}{k-1} \sum_{i=1}^{k} (Y_{ij} - \overline{Y}_{j})^{2} \\ V(\overline{Y}_{st}) &= \sum_{j=1}^{n} \frac{1}{n^{2}} \left(\frac{1}{1} - \frac{1}{k} \right) \frac{1}{k-1} \sum_{i=1}^{k} (Y_{ij} - \overline{Y}_{.j})^{2} \\ &= \sum_{i=1}^{k} \sum_{j=1}^{n} \frac{1}{n^{2}(k-1)} \left(\frac{k-1}{k} \right) (Y_{ij} - \overline{Y}_{.j})^{2} \\ &= \frac{1}{n^{2}k} \sum_{i=1}^{k} \sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{.j})^{2} \\ &= \frac{1}{nN} \sum_{i=1}^{k} \sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{.j})^{2} \\ n(k-1)V(\overline{Y}_{st}) &= \frac{n(k-1)}{nN} \sum_{i=1}^{k} \sum_{j=1}^{n} (Y_{ij} - \overline{Y}_{.j})^{2} \\ \end{split}$$

$$=S_{wst}^2 \quad \frac{N-n}{Nn}$$

On comparison we note that Systematic Sampling has the same precision as corresponding equivalent stratified sampling if $\rho_{wst} = 0$.

**** SPECIAL POPULATION:-**

If we have special population we have these type of situations.

The value of population increase accordingly to linear law i.e. differ by constant amount h(say).

Example:-Suppose t values of X are μ +h, μ +2h,...., μ +lh,..., μ +th Mean of this series is,

$$\begin{aligned} \text{Mean} &= \frac{1}{t} \sum_{i=1}^{t} (\mu + ih) \\ &= \frac{1}{t} \left[\mu t + h \frac{t(t+1)}{2} \right] \qquad \left| \sum_{i=1}^{t} i = \frac{t(t+1)}{2} \right] \\ \text{Mean} &= \mu + h \frac{(t+1)}{2} \\ \text{Mean Square} &= \frac{1}{N-1} \sum_{i} (y_i - \overline{Y})^2 \\ &= \frac{1}{t-1} \left[\sum_{i=1}^{t} \left\{ \mu + ih - \mu - h \frac{(t+1)}{2} \right\}^2 \right] \\ &= \frac{h^2}{t-1} \left[\sum_{i=1}^{t} \left\{ i^2 + i(t+1) + \frac{(t+1)^2}{4} \right\} \right] \\ &= \frac{h^2}{t-1} \left[\sum_{i} i^2 - (t+1) \sum_{i=1}^{t} i + \frac{t(t+1)^2}{4} \right] \\ &= \frac{h^2}{t-1} \left[\frac{t(t+1)(2t+1)}{6} - \frac{(t+1)t(t+1)}{2} + \frac{t(t+1)^2}{4} \right] \end{aligned}$$

$$= \frac{h^2}{t-1} \frac{t(t+1)}{12} [2(2t+1) - (t+1)6 + 3(t+1)]$$

$$= \frac{h^2}{t-1} \frac{t(t+1)}{12} [4t+2 - 6t - 6 + 3t + 3)]$$

$$= \frac{h^2}{12} \frac{t(t+1)}{t-1} [t-1]$$

$$= \frac{h^2}{12} t [t-1]$$

Theorem:-For special population with usual notation prove that $V_{st}:V_{sy}:V_{ran}::1/n:1:n \text{ if } k \text{ is large and hence,}$ $E_{st}:E_{sys}:E_{ran}::n^2:n:1$ **Proof:-(1) In SPSWOP, we have

Proof:-(1) In SRSWOR, we have,
$$V = -\frac{N-n}{S^2}$$

$$V_{ran} = \frac{N - n}{nN} S_y^2$$

= $\frac{N - n}{nN} \frac{1}{N - 1} \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \overline{Y}..)^2$

There are nk = N values in the population differing by constant amount

h.

$$V_{ran} = \frac{nk - n}{n^2 k} \frac{h^2}{12} nk(nk + 1)$$
$$= \frac{h^2(k - 1)(nk + 1)}{12}$$

(2) In systematic sampling, we have,

$$V_{sys} = \frac{1}{k} \sum_{i=1}^{k} (\bar{y}_i - \bar{Y}..)^2$$
$$= \frac{k-1}{k} \frac{1}{k-1} \sum_{i} (\bar{y}_i - \bar{Y}..)^2$$

 \overline{Y}_{i} differs by a constant amount h and there are such k values.

$$V_{sys} = \frac{k-1}{k} \frac{h^2 k(k-1)}{12}$$
$$= \frac{h^2 (k^2 - 1)}{12}$$

(3) In stratified sampling, we have.

$$V_{st} = \sum_{j=1}^{n} p_j^2 \left(\frac{1}{n_j} - \frac{1}{N_j} \right) S_j^2$$
$$= \sum_{j=1}^{n} p_j^2 \left(\frac{1}{n_j} - \frac{1}{N_j} \right) \frac{1}{k - 1} \sum_{i=1}^{k} (Y_{ij} - Y_{\cdot j})^2$$

There are k values each is differing by a constant h.

$$=\frac{\sum_{i=1}^{k}(Y_{ij}-Y_{j})^{2}}{k-1}=\frac{h^{2} k (k+1)}{12}$$

There are N = nk values divided in n strata each of size k and one unit per stratum to selected.

$$V_{st} = \sum_{j=1}^{n} \frac{1}{n^2} \left(\frac{1}{1} - \frac{1}{k} \right) \frac{h^2 k(k-1)}{12}$$
$$= \frac{h^2 (k^2 - 1)}{12n}$$

There fore

$$V_{st} : V_{sy} : V_{ran} :: \frac{h^{2}(k^{2}-1)}{12n} : \frac{h^{2}(k^{2}-1)}{12} : \frac{h^{2}(k-1)(nk+1)}{12n}$$

:: $\frac{k+1}{n} : (k+1) : nk+1$
: $\frac{1}{n} : 1 : n$ [coordinate of k is taken]
if k is large.

Next, $\frac{V_{ran}}{V_{st}} : \frac{V_{ran}}{V_{sy}} : \frac{V_{ran}}{V_{ran}} :: n^2 : n : 1$ if n is large, i.e. $E_{st} : E_{sy} : E_{ran} :: n^2 : n : 1$.

**** CIRCULAR SYSTEM SAMPLING:-**

The case N = nk is called linear systematic sampling. Let as examine what happens when $N \neq nk$.

e.g. N = 1	1, n = 3, k	= N/n = 11/3	$\cong 4$	
Random strata	System	Probability		
	i, i+k, i			
1	\mathbf{Y}_1	Y_5	Y9	1⁄4
2	\mathbf{Y}_2	Y_6	Y_{10}	1⁄4
3	Y_3	Y_7	Y_{11}	1⁄4
4	Y_4	Y_8		1⁄4

$$E(\overline{Y}_{sy}) = \frac{1}{k} \sum_{j=1}^{n} Y_i$$
$$= \frac{1}{4} \left[\overline{Y}_1 + \overline{Y}_2 + \overline{Y}_3 + \overline{Y}_4 \right]$$

$$\begin{split} &= \frac{1}{4} \left[\frac{Y_1 + Y_5 + Y_9}{3} + \frac{Y_2 + Y_6 + Y_{10}}{3} + \frac{Y_3 + Y_7 + Y_{11}}{3} + \frac{Y_4 + Y_8}{2} + \frac{Y_4 + Y_8}{3} - \frac{Y_4 + Y_8}{3} \right] \\ &= \frac{1}{4} \left[\sum_{i=1}^{11} \frac{Y_1}{3} + \left\{ Y_4 + Y_8 \right\} \left\{ \frac{1}{2} - \frac{1}{3} \right\} \right] \\ &= \frac{\sum_i Y_i}{12} + \frac{Y_4 + Y_8}{24} \\ &= \frac{11\overline{Y}_{..}}{12} + \frac{Y_4 + Y_8}{24} \\ &= \left(1 - \frac{1}{12} \right) \overline{Y}_{..} + \frac{Y_4 + Y_8}{24} \\ &= \overline{Y}_{..} + \left(\frac{Y_4 + Y_8}{24} - \frac{\overline{Y}_{..}}{12} \right) \\ E(\overline{Y}_{sys}) \neq \overline{Y}_{..} \end{split}$$

**** DRAW BACK OF LINEAR SYS SAMPLING (L.S.S.)**

- (1) All the sys sampling are not at the same size.
- (2) Sys sampling mean is not an unbiased estimator of the population mean.

To remove these draw back a rular Sys. Sampling is suggested. In C.S.S. the 1^{st} unit is selected by drawing a random no. from 1 to N and then select every kth unit in a cyclic manner till a sample of size n is obtain where,

K = the nearest integer N/n

i.e. $k \cong N/n$

e.g. N = 11, n = 3, $k = N/n \cong 4$

Suppose the first random no. is 4 then C.S.S. of size 3. Contains the unit as $Y_4,\,Y_8$, Y_i .

**** TWO STAGE SAMPLING ****

In cluster sampling the whole population is divided in to clusters and a sample of clusters is selected and the elements of the selected clusters are enumerated. If the elements of the clusters are more or less homogeneous then it is uneconomical to enumerate all the elements of the selected dusters. In this case two-stage or sub-sampling is suggested.

The sampling procedure of first selecting a sample of clusters and then selecting a sample of clusters and then selecting a sample of element from each selected clusters is known as two-stage or sub-sampling. The clusters which form sampling units at the first stage are called first-stage units and the elements within the cluster which form sampling units at the second stage are called second stage units.

For an example, while sampling of fields in a taluka, first a sample of villages is selected and from each selected village a sample of fields is selected.

This procedure can be extended and we have three stage or multi-stage sampling.

****** Notations :-

N = number of first stage units in the population.

M = number of second stage units in the population.

 Y_{ij} = Value of the jth second stage unit in the ith first stage unit.

$$i = 1, 2, ..., N.$$

$$j = 1, 2, ..., M.$$

$$\overline{Y}_{i} = \frac{1}{M} \sum_{j=1}^{M} Y_{ij}$$
 = Mean per second stage unit in ith first stage unit in the

population.

$$\overline{Y}_{..} = \frac{1}{NM} \sum_{i=1}^{N} \sum_{j=1}^{M} Y_{ij}$$
$$= \frac{1}{N} \sum_{i=1}^{N} \overline{Y}_{i}.$$

= Population mean = (mean per second stage unit in the whole population.)

$$S_i^2 = \frac{1}{M-1} \sum_{j=1}^M (Y_{ij} - \overline{Y}_i)^2$$

= Mean square among the second stage units in the i^{th} first age units in the population.

$$\overline{S}_{w}^{2} = \frac{1}{N} \sum_{i=1}^{N} S_{i}^{2} = \frac{1}{N(M-1)} \sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{i})^{2}$$

= Mean square among the second stage units within the first age units in the population.

$$S_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\overline{Y}_i - \overline{Y}_{..})^2$$

= Mean square among the first age units in the population.

We assume that a SRSWOR of n first age unit is known from the population and from each First stage unit a SRSWOR of m second stage unit is taken, thus we have a sample of nm second stage units.

$$\overline{Y}_{im} = \frac{1}{m} \sum_{j=1}^{m} Y_{ij}$$

= Mean per second stage unit in the i^{th} first age unit in the sample.

$$\overline{Y}_{nm} = \frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} Y_{ij} = \frac{1}{n} \sum_{i=1}^{n} \overline{Y}_{nm}$$

= Two stage sample mean. or mean per second stage unit in the i^{th} first stage unit in the sample.

$$s_i^2 = \frac{1}{m-1} \sum_{j=1}^m (y_{ij} - \overline{y}_{im})^2$$

= Mean square among the second stage units in the i^{th} first age units in the sample.

$$\int_{w}^{2} = \frac{1}{n} \sum_{i=1}^{n} s_{i}^{2} = \frac{1}{n(m-1)} \sum_{i=1}^{n} \sum_{j=1}^{m} (y_{ij} - \overline{y}_{im})^{2}$$

= Mean square among the second stage units within the i^{th} first age units in the sample.

$$s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{y}_{im} - \bar{y}_{nm})^2$$

= Mean square among the first stage units in the sample.

We know that
$$V(X) = EV(X/Y) + VE(X/Y)$$

 $E(X) = E_1 E_2(X/Y)$

* Theorem:- Prove that,

$$E(\overline{Y}_{nm}) = \overline{Y}..$$

$$V(\overline{Y}_{nm}) = \left(\frac{1}{n} - \frac{1}{N}\right)S_b^2 + \frac{1}{n}\left(\frac{1}{m} - \frac{1}{M}\right)\overline{S}_w^2$$

* Proof:-

$$\begin{split} E(\bar{Y}_{nm}) &= E\left[\frac{1}{n}\sum_{i=1}^{n}\bar{Y}_{im}\right] \\ &= E_{1}E_{2}\left[\frac{1}{n}\sum_{i=1}^{n}\bar{Y}_{im}/i\right] \\ &= E_{1}\left[\frac{1}{n}\sum_{i=1}^{n}E_{2}(\bar{Y}_{im}/i)\right] \\ &= E_{1}\left[\frac{1}{n}\sum_{i=1}^{n}F_{2}(\bar{Y}_{im}/i)\right] \\ &= E_{1}\left[\frac{1}{n}\sum_{i=1}^{n}\bar{Y}_{i}\right] \\ &= \frac{1}{n}\sum_{i=1}^{n}E_{1}[\bar{Y}_{i}.] \\ &= \frac{1}{n}\sum_{i=1}^{n}F_{1}. \\ &= \bar{Y}_{.}. \end{split}$$
(1)
$$V(\bar{Y}_{nm}) &= V_{1}E_{2}\left[\frac{1}{n}\sum_{i=1}^{n}\bar{Y}_{im}/i\right] + E_{1}V_{2}\left[\frac{1}{n}\sum_{i=1}^{n}\bar{Y}_{im}/i\right] \\ &= V_{1}\left[\frac{1}{n}\sum_{i=1}^{n}F_{2}(\bar{Y}_{im}/i)\right] + E_{1}\left[\frac{1}{n^{2}}\sum_{i=1}^{n}V_{2}(\bar{Y}_{im}/i)\right] \\ &= V_{1}\left[\frac{1}{n}\sum_{i=1}^{n}\bar{Y}_{.}\right] + \frac{1}{n^{2}}E_{1}\left[\sum_{i=1}^{n}\left(\frac{1}{m}-\frac{1}{M}\right)S_{i}^{2}\right] \\ &= \left(\frac{1}{n}-\frac{1}{N}\right)S_{b}^{2} + \frac{1}{n}\left(\frac{1}{m}-\frac{1}{M}\right)E_{1}\left[\frac{1}{n}\sum_{i=1}^{n}S_{i}^{2}\right] \\ &= \left(\frac{1}{n}-\frac{1}{N}\right)S_{b}^{2} + \frac{1}{n}\left(\frac{1}{m}-\frac{1}{M}\right)\overline{S}_{w}^{2} \end{aligned}$$
(ii)
So from (i) and (ii)
$$E(\bar{Y}_{nm}) = \left(\frac{1}{n}-\frac{1}{N}\right)S_{b}^{2} + \frac{1}{n}\left(\frac{1}{m}-\frac{1}{M}\right)\overline{S}_{w}^{2} \\ Estimation of V(\bar{Y}_{nm}) \end{split}$$

Theorem:- An unbiased estimator of the variance of two-stage sample mean is given by,

$$Est \ V(\overline{Y}_{nm}) = \left(\frac{1}{n} - \frac{1}{N}\right) s_b^2 + \frac{1}{N} \left(\frac{1}{m} - \frac{1}{M}\right) \overline{s}_w^2$$

where, $s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\overline{y}_{im} - \overline{y}_{nm})^2$
 $s_i^2 = \frac{1}{m-1} \sum_{j=1}^m (y_{ij} - \overline{y}_{im})^2$
 $\overline{s}_w^2 = \frac{1}{n} \sum_{i=1}^n s_i^2$

Proof:-

$$V(\overline{V}_{nm}) = \left(\frac{1}{n} - \frac{1}{N}\right) s_b^2 + \frac{1}{n} \left(\frac{1}{m} - \frac{1}{M}\right) \overline{s}_w^2$$

i.e. Est $V(\overline{Y}_{nm}) = \left(\frac{1}{n} - \frac{1}{N}\right) Est s_b^2 + \frac{1}{n} \left(\frac{1}{m} - \frac{1}{M}\right) Est \overline{s}_w^2$ (iii)
Now, $E(\overline{s}_w^2) = E\left[\frac{1}{n}\sum_{i=1}^n s_i^2\right]$
 $= E_1 E_2 \left[\frac{1}{n}\sum_{i=1}^n S_i^2/i\right]$
 $= E_1 \left[\frac{1}{n}\sum_{i=1}^n E_2(s_i^2/i)\right]$
 $= E_1 \left[\frac{1}{n}\sum_{i=1}^n S_i^2\right]$
 $= \frac{1}{N}\sum_{i=1}^N S_i^2 = \overline{S}_w^2$ (iv)
and $s_b^2 = \frac{1}{n-1}\sum_{i=1}^n (\overline{y}_{im} - \overline{y}_{nm})^2$
i.e. $(n-1) E(s_b^2) = E\left[\sum_{i=1}^n (\overline{y}_{im})^2\right] - nE\left[\sum_{i=1}^n (\overline{y}_{nm}) + \overline{Y}^2.\right]$
[Since $V(\overline{Y}_{nm}) = E(\overline{Y}_{nm}^2) - [E(\overline{Y}_{nm})^2] = E(\overline{Y}_{nm}^2) - \overline{Y}.^2$
So, $E(\overline{Y}_{nm}^2) = V(\overline{Y}_{nm}) + \overline{Y}.^2$]
 $= E_1 \left[\sum_{i=1}^n E_2(\overline{y}_{im}^2/i)\right] - n\left[\left(\frac{1}{n} - \frac{1}{N}\right)S_b^2 + \frac{1}{n}\left(\frac{1}{m} - \frac{1}{M}\right)\overline{S}_w^2\right] - n\overline{Y}.^2$

$$\begin{split} &= E_{1} \left[\sum_{i=1}^{n} \left\{ V_{2}(\bar{y}_{im}^{2}/i) + \left(E_{2}(\bar{y}_{im}/i) \right)^{2} \right\} \right] - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} - \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \bar{Y}..^{2} \\ &= E_{1} \left[\sum_{i=1}^{n} \left\{ \left(\frac{1}{m} - \frac{1}{M} \right) S_{i}^{2} + \bar{Y}_{i}.^{2} \right\} \right] - n \bar{Y}..^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} - \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} \\ &= n \left(\frac{1}{m} - \frac{1}{M} \right) E_{1} \left[\frac{1}{n} \sum_{i=1}^{n} S_{i}^{2} \right] + n E_{1} \left[\frac{1}{n} \sum_{i} \bar{Y}_{i}.^{2} \right] - n \bar{Y}..^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} - \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} \\ &= n \left(\frac{1}{m} - \frac{1}{M} \right) \left[\frac{1}{N} \sum_{i=1}^{N} S_{i}^{2} \right] + n E_{1} \left[\frac{1}{n} \sum_{i} \bar{Y}_{i}.^{2} \right] - n \bar{Y}..^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} - \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} \\ &= n \left(\frac{1}{m} - \frac{1}{M} \right) \left[\frac{1}{N} \sum_{i=1}^{N} S_{i}^{2} \right] + n E_{1} \left[\frac{1}{n} \sum_{i} \bar{Y}_{i}.^{2} \right] - n \bar{Y}..^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} - \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} \\ &= n \left(\frac{1}{m} - \frac{1}{M} \right) \left[\bar{X}_{w}^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} - \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} \\ &= \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} + \frac{n}{N} \left(\sum_{i=1}^{N} \bar{Y}_{i}.^{2} - N \bar{Y}..^{2} \right) \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} + \frac{n}{N} \left(N - 1 \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \left(\frac{1}{n} - \frac{1}{N} \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \left(\frac{n - 1}{n} \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - n \left(\frac{n - 1}{n} \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - \left(n - 1 \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - \left(n - 1 \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - \left(n - 1 \right) S_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - \left(n - 1 \right) \bar{S}_{b}^{2} \\ &= \left(n - 1 \right) \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_{w}^{2} - \left(n - 1 \right$$

** Efficiency of Two-Stage Sampling and SRSWOR:-

In SRSWOR we have,

$$V(\overline{Y}_{nm}) \ SRSWOR = \left(\frac{1}{nm} - \frac{1}{NM}\right)S^2$$
$$\cong \ \frac{1}{nm}S^2 \ (\frac{1}{M} \text{ is negligible}) \qquad \text{or (M is large.)}$$

In two-stage sampling, we have,

$$V(\overline{Y}_{nm}) \ two - stage = \left(\frac{1}{n} - \frac{1}{N}\right)S_b^2 + \frac{1}{n}\left(\frac{1}{m} - \frac{1}{M}\right)\overline{S}_w^2$$
$$\cong \left(\frac{1}{n} - \frac{1}{N}\right)S_b^2 + \frac{1}{nm}\overline{S}_w^2 \quad (\frac{1}{M} \ is \ negligible)$$
we have M^2 (N-1) $S_b^2 = (NM-1)S_b^2 [1+(M-1)\alpha]$

now we have M^2 (N-1) $S_b^2 = (NM-1) S^2 [1+(M-1)\rho]$

i.e. $S_b^2 = \frac{NM - 1}{M(N - 1)} S^2 \left[\frac{1 + (M - 1)\rho}{M} \right]$ $= \frac{N}{N - 1} S^2 \rho \qquad (\frac{1}{M} \text{ is negligible})$

Next we have,

$$(NM - 1) S^{2} = N(M - 1) \overline{S}_{w}^{2} + M(N - 1) S_{b}^{2}$$

$$\overline{S}_{w}^{2} = \frac{(NM - 1) S^{2} - M(N - 1) S_{b}^{2}}{N(M - 1)}$$

$$= \frac{\left(N - \frac{1}{M}\right) S^{2} - (N - 1) S_{b}^{2}}{N\left(1 - \frac{1}{M}\right)}$$

$$= \frac{NS^{2} - (N - 1) S_{b}^{2}}{N}$$

$$= S^{2} - \left(\frac{N - 1}{N}\right) S_{b}^{2} \qquad (\because \frac{1}{M} \text{ is negligible })$$

$$= S^{2} - \frac{N - 1}{N} \frac{N}{N - 1} S^{2} \rho$$

$$= S^{2} - S^{2} \rho$$

$$= S^{2} (1 - \rho)$$

There fore $V(\overline{Y}_{nm})$ two-stage becomes.

$$V(\overline{Y}_{nm}) \text{ two stage} \cong \frac{N-n}{nN} \frac{N}{N-1} S^2 \rho + \frac{S^2}{nm} (1-\rho)$$
$$\cong \frac{N-n}{n(N-1)} S^2 \rho + \frac{S^2}{nm} (1-\rho)$$
$$\cong \frac{S^2}{nm} \left[\frac{m(N-n)}{N-1} \rho + 1 - \rho \right]$$
$$\cong \frac{S^2}{nm} \left[1 + \left\{ \frac{m(N-n)}{N-1} - 1 \right\} \rho \right]$$

Efficiency of two-stage sampling over SRSWOR is given by,

Efficiency
$$\approx \frac{V(Y_{nm}) \ SRSWOR}{V(\overline{Y}_{nm}) \ two-stage}$$

$$\approx \frac{S^2/nm}{\frac{S^2}{nm} \left[1 + \left\{\frac{m(N-n)}{N-1} - 1\right\}\rho\right]}$$

Eff.
$$\cong \frac{1}{1 + \left\{\frac{m(N-n)}{N-1} - 1\right\}\rho}$$

Obtain the efficiency of two stage sampling compared to SRSWOR in forms of intra class correlation to efficient.

Optimum Value of n and m :- (For Two-Stage Sampling)

$$V(\overline{Y}_{nm}) = \left(\frac{1}{n} - \frac{1}{N}\right) s_b^2 + \frac{1}{n} \left(\frac{1}{m} - \frac{1}{M}\right) \overline{s}_w^2$$

Consider the cost function.

 $\mathbf{c} = \mathbf{c}_1 \mathbf{n} + \mathbf{c}_2 \mathbf{n} \mathbf{m}.$

Where c = total cost

 $c_1 = cost$ for collecting information per first stage unit.

 $c_2 = cost$ for collecting information per second stage unit. Determine the optimum value of n and m such that variance.

i.e $V(\overline{Y}_{nm})$ is minimized for fined cost w, say,

$$\phi = \left(\frac{1}{n} - \frac{1}{N}\right) s_b^2 + \frac{1}{n} \left(\frac{1}{m} - \frac{1}{M}\right) \bar{s}_w^2 + \lambda [c_1 n + c_2 nm] = 0$$

$$\frac{\partial \phi}{\partial n} = -\frac{1}{n^2} \bar{s}_b^2 + \left(\frac{1}{m} - \frac{1}{M}\right) \bar{s}_w^2 + \lambda (c_1 + c_2 m) = 0$$

$$\Rightarrow n^2 = \frac{s_b^2 + \left(\frac{1}{m} - \frac{1}{M}\right) \bar{s}_w^2}{\lambda [c_1 + c_2 m]} \qquad (i)$$

$$\frac{\partial \phi}{\partial m} = -\frac{1}{nm^2} \bar{s}_w^2 + \lambda c_2 n = 0$$

$$\Rightarrow -\lambda c_2 n = -\frac{\bar{s}_w^2}{nm^2}$$

$$\Rightarrow n^2 = \frac{\bar{s}_w^2}{\lambda c_2 m^2} \qquad (ii)$$

From (i) and (ii) we have,

$$\frac{s_b^2 + \left(\frac{1}{m} - \frac{1}{M}\right)\overline{s}_w^2}{\lambda[c_1 + c_2mn]} = \frac{\overline{s}_w^2}{\lambda c_2m^2}$$

i.e $c_2m^2 \left[s_b^2 + \left(\frac{1}{m} - \frac{1}{M}\right)\overline{s}_w^2\right] = [c_1 + c_2m] \ \overline{s}_w^2$
i.e $c_2m^2 \left[s_b^2 - \left(\frac{1}{M} \ \overline{s}_w^2\right)\right] = c_1 \ \overline{s}_w^2$

i.e
$$c_2 m^2 D = c_1 \bar{s}_w^2$$
, where $D = s_b^2 - \frac{1}{M} \bar{s}_w^2$
 $m^2 = \frac{c_1 \bar{s}_w^2}{c_2 D}$
If $D > 0$
 $\hat{m} = \sqrt{\frac{c_1 \bar{s}_w^2}{c_2 D}}$

Now, $c_0 = c_1 n + c_2 n \hat{m}$

i.e.
$$\hat{n} = \frac{c_0}{c_1 + c_2 \hat{m}}$$

If D < 0, the above procedure fails in this case we proceed as follows.

$$V(\overline{Y}_{nm}) = \left(\frac{1}{n} - \frac{1}{N}\right) s_b^2 + \frac{1}{n} \left(\frac{1}{m} - \frac{1}{M}\right) \overline{s}_w^2$$
$$= \frac{1}{n} \left[s_b^2 - \frac{1}{M} \overline{s}_w^2\right] + \frac{1}{nm} \overline{s}_w^2 - \frac{1}{N} s_b^2$$
$$= \frac{1}{n} D + \frac{1}{nm} \overline{s}_w^2 - \frac{1}{N} s_b^2$$

now, $c_0 = c_1 n + c_2 nm$.

i.e.
$$n = \frac{c_0}{c_1 + c_2 m}$$
(iii)

Substituting the above , value of n from (iii) becomes,

$$V(\overline{Y}_{nm}) = \frac{c_1 + c_2 m}{c_0} D + \frac{c_1 + c_2 m}{c_0 m} \overline{S}_w^2 - \frac{1}{N} S_b^2$$

As m increases, $V(\overline{Y}_{nm})$ decreases and $V(\overline{Y}_{nm})$ is minimum when m is the largest positive integer.

i.e. $\hat{m} = M$. Then, n is obtained from the following.

$$\hat{n} = \frac{c_0}{c_1 + c_2 \hat{m}}$$
$$= \frac{c_0}{c_1 + c_2 M}$$

So $\hat{n} = \frac{c_0}{c_1 + c_2 \hat{m}}$ and $\hat{m} = M$ optimum value of n and m.

**** CLUSTER SAMPLING****

A sample procedure pre-suppose the division of the population in to a finik no. of distinct and identifying units that called sampling unit. The smallest units in to which the population can be divided are called the element of population and the group of element are called clusters. Some times instead of taking an element as the sampling unit cluster is taken as the sampling unit. This is useful when the list of element is not available while the list of cluster is available.

Thus, in cluster sampling, A population is divided into cluster and sample random sample of clusters is selected and information for all elements of the selected cluster is collected.

EXAMPLE:-

While sampling of population is a certain city, the list of person residing in that city is not available, but the list of households is readily available.

Value of study variable y in the population:-

	1	2		i		Ν
1	Y_{11}	Y_{21}		Y_{i1}		Y_{N1}
2	Y_{12}	Y_{22}	••••	Y_{i2}	•••••	$Y_{N2} \\$
•	•	•	•	•	•	•
•	•	•	•	•	•	•
j	\mathbf{Y}_{1j}	Y_{2j}	••••	\mathbf{Y}_{ij}	•••••	Y_{Nj}
•	•	•	•	•	•	•
•	•	•	•	•	•	•
Μ	Y_{1M}	Y_{2M}		Y_{iM}		Y_{NM}

NM = Number of elements in the population.

N = Number of clusters in the population.

M = Number of elements in each clusters.

 Y_{ij} = Value of jth elements in the ith cluster (i= 1,2,...,N. j = 1,2,...,M.)

$$\overline{Y}_{i} = \frac{1}{M} \sum_{j=1}^{M} Y_{ij} = \text{Mean of the } i^{\text{th}} \text{ cluster.}$$

$$\overline{Y}_{NM} = \frac{1}{NM} \sum_{i=1}^{N} \sum_{j=1}^{M} Y_{ij} = \frac{1}{N} \sum_{i=1}^{N} \overline{Y}_{i}.$$

$$= \text{population mean.}$$

$$S_{i}^{2} = \frac{1}{M-1} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{NM})^{2}$$

= mean square between elements in i^{th} cluster.

$$S_w^2 = \frac{1}{N} \sum_{i=1}^N S_i^2$$

= mean square between elements within the clusters in the population.

$$S^{2} = \frac{1}{NM - 1} \sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{NM})$$

= Mean square between elements in the whole population.

$$S_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\overline{Y}_i - \overline{Y}_{NM})^2$$

= Mean square between clusters in the whole population.

Suppose a simple random sample of n cluster is drawn from N elements and all the elements of selected clusters are enumerated. Thus, we have a cluster sample of NM elements.

$$\overline{Y}_{nM} = \frac{1}{nM} \sum_{i} \sum_{j} y_{ij}$$
$$= \frac{1}{n} \sum_{i} \overline{Y}_{i} = \text{Sample mean.}$$
$$\overline{s}_{w}^{2} = \frac{1}{n} \sum_{i} s_{i}^{2}$$
$$= \frac{1}{n(M-1)} \sum_{i} \sum_{j} (Y_{ij} - \overline{Y}_{..})^{2}$$

= Mean square between elements within the cluster in the

sample.

$$s_b^2 = \frac{1}{n-1} \sum_i (y_{ij} - \overline{Y}_{nM})^2$$

= Mean square between clusters in the sample.
$$s^2 = \frac{1}{nM-1} \sum_i \sum_j (y_{ij} - \overline{Y}_{nM})^2$$

= Mean square between elements in the whole sample.

* Theorem:-Prove that-

1)
$$E(\overline{Y}_{nM}) = \overline{Y}_{NM}$$

2) $E(\overline{Y}_{nM}) = \left(\frac{1}{n} - \frac{1}{N}\right) S_b^2$

* **Proof:-**Let Q_i be the quantity associated with the ith unit of the population (i= 1,2,...,N)

Let $\overline{Q} = \frac{1}{N} \sum_{i=1}^{N} Q_i$ = population mean. $S_Q^2 = \frac{1}{N-1} \sum_{i=1}^{N} (Q_i - \overline{Q}_N)^2$

Suppose a simple random sample of size n is taken from the population of size N.

Let
$$\overline{Q}_n = \frac{1}{n} \sum_{i=1}^n Q_i$$

 $S_Q^2 = \frac{1}{n-1} \sum_{i=1}^n (Q_i - \overline{Q}_n)^2$

For **SRSWOR** we have,

i)
$$E(\overline{Q}_n) = \overline{Q}_N$$

ii) $V(\overline{Q}_n) = \left(\frac{1}{n} - \frac{1}{N}\right) S_Q^2$
iii) $E(s_Q^2) = S_Q^2$

For cluster sampling, let us define $Q_i = Y_i$, $\overline{Q}_n = \overline{Y}_{nM}$, $\overline{Q}_N = \overline{Y}_{NM}$, $S_Q^2 = S_b^2$

Appling (i) and (ii) we get,

$$E(\overline{Y}_{nM}) = \overline{Y}_{NM}$$
$$V(\overline{Y}_{nM}) = \left(\frac{1}{n} - \frac{1}{M}\right) S_b^2, \text{ Hence proved.}$$

** Efficiency of Cluster Sampling Compared To SRSWOR:-

In **SRSWOR**, we have

$$V(\overline{Y}_{nM}) = \left(\frac{1}{nm} - \frac{1}{NM}\right) S^2$$

For cluster sampling, we have

$$V(\overline{Y}_{nM}) = \left(\frac{1}{n} - \frac{1}{M}\right) S_b^2$$

Efficiency of cluster sampling over SRSWOR is given by:

$$E = \frac{V(\overline{Y}_{nM})}{V_{ci}(\overline{Y}_{nM})} = \frac{1/M\left(\frac{1}{n} - \frac{1}{N}\right)S^2}{\left(\frac{1}{n} - \frac{1}{N}\right)S_b^2} = \frac{S^2}{MS_b^2}$$

(1) We have,

$$(1) \text{ we have,}$$

$$(NM - 1) s^{2} = \sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{NM})^{2}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{i} + \overline{Y}_{i} - \overline{Y}_{NM})^{2}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{i})^{2} + \sum_{i} \sum_{j} (\overline{Y}_{i} - \overline{Y}_{NM})^{2} + \text{ product term}$$

$$(NM - 1)S^{2} = N(M - 1) S_{w}^{-2} + M(N - 1) S_{b}^{2}$$
where
$$\rho = \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{i})(Y_{ij} - \overline{Y}_{i})/NM(M - 1)}{\sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{NM}^{2})/NM}$$

$$= \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} (Y_{ij} - \overline{Y}_{NM})(Y_{ij}^{'} - \overline{Y}_{NM})}{(M - 1)(NM - 1)S^{2}}$$

= Intra class (intra cluster) correlation co efficient between pairs of elements within clusters.

From (i) and (ii) we get (NM-1) $S^2 [1+(M-1) \rho] = M^2 (N-1) S_b^2$ $\therefore E = \frac{S^2}{MS_b^2} = \frac{M(N-1)}{(NM-1)[1+(M-1)\rho]}$ if N is large $E \approx 1/1 + (M-1)\rho$ i.e. $S^2 = \frac{N(M-1)}{NM-1} S_w^{-2} + \frac{M(N-1)}{NM-1} S_b^2$ we have, Efficiency = $E = \frac{S^2}{MS_b^2}$ $= \frac{\frac{N(M-1)}{NM-1} S_w^{-2} + \frac{M(N-1)}{NM-1} S_b^2}{MS_b^2}$ $= \frac{N(M-1)}{M(NM-1)} \frac{S_w^{-2}}{S_b^2} + \frac{N-1}{NM-1}$ $E \approx \frac{M-1}{M} \frac{S_w^{-2}}{S_b^2} + \frac{1}{M}$; if N is large.

* ANOVA TABLE *

Source of	Sum of squares	Degree of	M.S.S.	Е
variation		fraction		
Between cluster	$\sum_{i}\sum_{j}(Y_{i}\overline{Y}_{NM})^{2}$	N-1	MS_b^2	S^2/MS_b^2
Within cluster	$\sum_{i}\sum_{j}(Y_{ij}-\overline{Y}_{i}.)^{2}$	N(M-1)	S_w^{-2}	
Total	$\sum_{i}\sum_{j}(Y_{ij}-\overline{Y}_{NM})^{2}$	NM-1	S^2	

(2) Efficiency of cluster sampling over SRSWOR in forms of Intra-class correlation co-efficient is given by:

$$E \cong \frac{1}{1 + (M+1)\rho}$$

**** ESTIMATION OF EFFICIENCY:-**

A estimate of the efficiency of cluster sampling over SRSWOR is given by:

$$\hat{E} = \frac{M-1}{M} \frac{s_w^{-2}}{MS_b^2} + \frac{1}{M}$$
 if N is large.

Proof:- We have to show that

$$E(s_w^{-2}) = S_w^{-2}$$
 and $E(s_b^2) = S_b^2$

Let Q_i be the quantity associated with the i^{th} unit of the population $i = 1, 2, \dots, N$.

Let
$$\overline{Q}_N = \frac{1}{N} \sum_{i=1}^N Q_i$$

 $S_Q^2 = \frac{1}{N-1} \sum_{i=1}^N (Q_i - \overline{Q}_n)^2$

Suppose a simple random sample of n unit is drawn from N unit.

Let
$$\overline{Q}_n = \frac{1}{n} \sum_{i=1}^n Q_i$$

 $s_Q^2 = \frac{1}{n-1} \sum_{i=1}^n (Q_i - \overline{Q}_n)^2$

For SRSWOR we have,

i)
$$E(\overline{Q}_n) = \overline{Q}_N$$

ii) $V(\overline{Q}_n) = \left(\frac{1}{n} - \frac{1}{N}\right) S_Q^2$
iii) $E(s_Q^2) = S_Q^2$

For cluster sampling let us define.

$$Q_i = S_i^2$$
, $\overline{Q}_n = \frac{1}{n} \sum_{i=1}^n S_i^2 = S_w^{-2}$, $\overline{Q}_N = S_w^2$,

Appling (i) we get,

$$E(s_w^{-2}) = S_w^{-2}$$

Again for cluster sampling let us define

 $Q_i = \overline{Y}_i, \ \overline{Q}_n = \overline{Y}_{nm}, \ \overline{Q}_N = \overline{Y}_{NM}, \ s_Q^2 = s_b^2, \ S_Q^2 = S_b^2$ Appling (i) we have $E(s_b^2) = S_b^2$

**** OPTIMUM VALUE OF n AND M:-**

Faifield Smith (1938).

 $S_b^2 = \frac{S^2}{M}$ is some constant less than 1.

Mahalnibis (1940) and Jessen (1942)

 $S_w^{-2} = a.M^b$ a and b are constant greater then zero.

We have
$$\sum_{i=1}^{N} \sum_{j=1}^{M} (y_{ij} - \overline{Y}_{NM})^2 = \sum_{i} \sum_{j=1}^{N} (y_{ij} - \overline{Y}_{i})^2 + \sum_{i} \sum_{j=1}^{N} (\overline{Y}_{i} - \overline{Y}_{NM})^2$$

i.e.
$$(NM-1) S^{2} = N(M-1) S_{w}^{-2} + M(N-1) S_{b}^{2}$$

 $S_{b}^{2} = \frac{(NM-1)S^{2}N(M-1)S_{w}^{-2}}{M(N-1)}$
 $\cong S^{2} - S_{w}^{-2} \frac{(M-1)}{M}$ if N is large
 $\cong S^{2} - a(M-1)(M)^{b-1}$; $S_{w}^{-2} = a.M^{b}$
 $V(\overline{Y}_{nM}) = \left(\frac{1}{n} - \frac{1}{N}\right) S_{b}^{2}$
 $\cong \frac{1}{n} S_{b}^{2}$ for large N
 $\cong \frac{1}{n} [S^{2} - a(M-1)M^{b-1}] \cong V$

let the cost function bc

 $c = c_1 n M + c_2 \sqrt{n}$

where, $c_1 = Cost$ for colleting information per element

 $c_2 = Cost proportional to unit distance between cluster.$

Our problem is to determine the optimum value of n and M, so that variance is minimized for fixed cost c_0 say.

We minimize,

$$\sqrt{n} = \frac{-c_2 \pm \sqrt{\Delta}}{2aM} \quad , \qquad \Delta = c_2^2 + 4c_1 c_0 M$$

$$= \frac{-c_{2} + \sqrt{\Delta}}{2aM} \qquad (\because \sqrt{n} \text{ being positive})$$

now, $c_{1}M + \frac{c_{2}}{2\sqrt{n}} = c_{1}M + \frac{c_{2} 2c_{1}M}{2(-c_{2} + \sqrt{\Delta})}$
 $= c_{1}M + \frac{c_{2}c_{1}M}{-c_{2} + \sqrt{\Delta}}$
 $= \frac{c_{1}M\sqrt{\Delta}}{-c_{2} + \sqrt{\Delta}}$
(2) reduce to i: $\frac{1}{V} \frac{\partial V}{\partial M} = \frac{-c_{1}[-c_{2} + \sqrt{\Delta}]}{aM\sqrt{\Delta}}$
 $\Rightarrow \frac{M}{V} \frac{\partial V}{\partial M} = \frac{-c_{2} + \sqrt{\Delta}}{\sqrt{\Delta}}$
 $\Rightarrow \frac{M}{V} \frac{\partial V}{\partial M} = \frac{c_{2} - 1}{\sqrt{\Delta}}$
 $\Rightarrow \frac{M}{V} \frac{\partial V}{\partial M} = \frac{c_{2}}{\sqrt{c_{2}^{2} + 4c_{1}c_{0}M}} - 1$
 $\Rightarrow \frac{M}{V} \frac{\partial V}{\partial M} = \frac{c_{2}}{\sqrt{c_{2}^{2}(\frac{1 + 4c_{1}c_{0}M}{c_{2}^{2}})}} - 1$
 $= \frac{1}{\sqrt{\frac{1 + 4c_{1}c_{0}M}{c_{2}^{2}}}} - 1$
 $\therefore \frac{M}{V} \frac{\partial V}{\partial M} = -1 + \left[\frac{1 + 4c_{1}c_{0}M}{c_{2}^{2}}\right]^{-1/2} \qquad \dots (3)$

$$\frac{\partial V}{\partial M} = \frac{\partial}{\partial M} \left[\frac{1}{n} \left\{ S^2 - a(M-1)M^{b-1} \right\} \right]$$
$$= \frac{1}{n} \left[0 - \left\{ aM^{b-1} + a(M-1)(b-1)M^{b-2} \right\} \right]$$
$$= \frac{1}{n} \left[-aM^{b-1} - a(M-1)(b-1)M^{b-2} \right]$$
$$\Rightarrow M \frac{\partial V}{\partial M} = \frac{1}{n} \left[-aM^b - a(M-1)(b-1)M^{b-1} \right]$$
$$\Rightarrow \frac{M}{v} \frac{\partial V}{\partial M} = \frac{1/n \left[-aM^b - a(M-1)(b-1)M^{b-1} \right]}{1/n \left[S^2 - a(M-1)M^{b-1} \right]}$$

$$=\frac{\left[-aM^{b}-a(M-1)(b-1)M^{b-1}\right]}{\left[S^{2}-a(M-1)M^{b-1}\right]}$$
(4)

From (3) and (4) we find

$$-1 + \left[\frac{1 + 4c_1c_0M}{c_2^2}\right]^{-1/2} = \frac{-aM^b - a(M-1)(b-1)M^{b-1}}{S^2 - a(M-1)M^{b-1}}$$

This equation is to be solved by mal and error method. It will give the optimum value of M. Optimum value M is called \hat{M} say.

Then the optimum value of n, \hat{n} say is obtained by,

$$\sqrt{n} = \frac{-c_2 + \sqrt{c_2^2 + 4c_1c_0\hat{M}}}{2a\hat{M}}$$

For cluster sampling obtain for minimize variable cost is fixed.

** DOUBLE SAMPLING ** OR ** SUB SAMPLING **

Question:-Why double sampling is used or necessary ?

Answer:-When the information auxiliary variate is available, we have seen how it could be utilized to obtain the more efficient estimator.

e.g. Ratio and regression methods estimators.

If the information on auxiliary variable x is not available, we select an initial sample from the population and information on x is collected. Then we consider a second sample from the initial sample and collect the information on y. This procedure of first selecting an initial sample and then a second sample from the initial sample is known as double sampling or sub-sampling.

* Double Sampling for PPS Sampling. :-

Suppose we wish to estimate the population mean \overline{Y} using probability proportioned to size with replacement (PPSWR) scheme. Where the size being the value of some auxiliary variate x. Suppose the information on auxiliary variate x is not available. Then, we select an initial sample of size n' using SRSWOR and collect the information on x. Then, we select a second sample of size n from the initial sample using probability proportional to size with replacement scheme.

Then,
$$\hat{\overline{Y}} = \frac{1}{n} \sum z_i$$
, $z_i = \frac{Y_i}{np_i}$, $p_i = \frac{x_i}{x'}$

Is an unbiased estimators of \overline{Y} .. with variance is given that,

 S_y^2

$$Var (\hat{\overline{Y}}) = \frac{1}{N(N-1)} V_p(Y) \frac{n'-1}{nn'} + \frac{N-n'}{n'N}$$

where $V_P(Y) = \sum_{i=1}^{N} \frac{x_i}{X} \left(\frac{Y_i}{x_i/X} - Y\right)^2$
 $S_y^2 = \frac{1}{N-1} \sum_{i=1}^{N} (Y_i - \overline{Y})^2$
 $x' = \sum_{i=1}^{n'} x_i$
 $Y = \sum_{i=1}^{N} Y_i$
 $X = \sum_{i=1}^{N} X_i$
Now, $E(\hat{\overline{Y}}) = E\left[\frac{1}{n} \sum_{i=1}^{n} z_i\right]$
 $= E_1 E_2\left[\frac{1}{n} \sum_{i=1}^{n} z_i\right]$

 $E_2 \Rightarrow$ second sample.

 $E_1 \Rightarrow initial (first) sample.$

When E_2 denotes the conditional expectation over second sample when initial sample is fixed.

$$E(\hat{Y}) = E_{1}(\hat{y}) \qquad E_{1}E_{2}\left[\frac{1}{n}\sum_{i=1}^{n}z_{i}\right]$$

$$= \overline{Y} \qquad E_{1}[E_{2} \ z_{i}] = E_{1}[\overline{y}'] = \overline{Y}$$
Next $Var(\hat{Y}) = E_{1}V_{2}(\hat{y}) + V_{1}E_{2}(\hat{y})$
Consider $V_{1}E_{2}[\hat{y}] = V_{1}[\overline{Y}']$

$$= \frac{N-n'}{n'N}S_{y}^{2} \qquad (i)$$

Further the conditional variance of $\hat{\vec{Y}}$ when the initial sample fixed is given by,

$$V_{z}(\hat{Y}) = \frac{1}{n} \sum_{i=1}^{n'} \frac{x_{i}}{x'} \left(\frac{Y_{i}}{n'x_{i}/x'} - Y^{-1} \right)^{2}$$

$$= \frac{1}{nn'^{2}} \sum_{i=1}^{n'} \frac{x_{i}}{x'} \left(\frac{Y_{i}x'}{x_{i}} - Y' \right)^{2} \qquad (\because n'Y^{-1} = Y')$$
$$= \frac{1}{nn'^{2}} \sum_{i=1}^{n'} x_{i}x' \left(\frac{Y_{i}}{x_{i}} - \frac{Y'}{x'} \right)^{2}$$
$$= \frac{1}{nn'^{2}} \sum_{i=1}^{n'} \sum_{i< j}^{n'} x_{i}x_{j} \left(\frac{Y_{i}}{x_{i}} - \frac{Y_{j}}{x_{j}} \right)^{2}$$
$$= \frac{1}{nn'^{2}} \sum_{i=1}^{N} \sum_{i< j}^{N} a_{i}a_{j} x_{i}x_{j} \left(\frac{Y_{i}}{x_{i}} - \frac{Y_{j}}{x_{j}} \right)^{2}$$

Where $a_j = 1$ if the ith unit is included in the sample. $a_i = 0$ if the ith unit is not included in the sample.

$$E(a_i) = \frac{n'}{N}$$
$$E(a_i, a_j) = \frac{n'(n'-1)}{N(N-1)}$$

for a SRSWOR of size n'

$$E_{1}V_{2}[\hat{Y}] = 1 \left[\frac{1}{nn'^{2}} \sum_{i=1}^{N} \sum_{i\neq j}^{N} a_{i}a_{j} x_{i}x_{j} \left(\frac{y_{i}}{x_{i}} - \frac{y_{j}}{x_{j}} \right)^{2} \right]$$

$$= \frac{1}{nn'^{2}} \frac{n'(n'-1)}{N(N-1)} \sum_{i=1}^{N} \sum_{i\neq j}^{N} x_{i}x_{j} \left(\frac{y_{i}}{x_{i}} - \frac{y_{j}}{x_{j}} \right)^{2}$$

$$= \frac{1}{N(N-1)} \frac{n'-1}{nn'} \sum_{i=1}^{N} x_{i}X \left(\frac{y_{i}}{x_{i}} - \frac{Y}{X} \right)^{2}$$

$$= \frac{1}{N(N-1)} \frac{n'-1}{nn'} \sum_{i=1}^{N} \frac{x_{i}}{X} \left(\frac{y_{i}X}{x_{i}} - Y \right)^{2}$$

$$= \frac{1}{N(N-1)} \frac{n'-1}{nn'} \sum_{i=1}^{N} \frac{x_{i}}{X} \left(\frac{y_{i}}{x_{i}/X} - Y \right)^{2}$$

$$= \frac{1}{N(N-1)} \frac{n'-1}{nn'} V_{p}(Y) \qquad (ii)$$

Substitution (i) and (ii)

Var $(\hat{\overline{Y}})$ becomes;

$$Var \ (\hat{\overline{Y}}) = \frac{1}{N(N-1)} \frac{n'-1}{nn'} V_p(Y) + \frac{N-n'}{n'N} S_y^2$$

**** DOUBLE SAMPLING FOR UNBIASED RATIO ESTIMATOR:-**

Suppose we wish to estimate the population mean \overline{Y} using ratio method of estimation. Suppose the information on auxiliary variable x is not available.

Select a SRSWOR of size n' from the population and collect the information on y and \boldsymbol{x}

Consider the estimator as $\hat{Y} = \frac{\bar{Y}}{\bar{x}} x'$

Where $\overline{Y}(\overline{x})$ = mean of the second sample of size n of variate Y(x)

x' = mean of the initial sample of size n' of variate x.

now, $E(\hat{\overline{Y}}) = E\left[\frac{\overline{Y}}{\overline{x}} \ \overline{x}'\right]$ $= E_1 E_2 \left[\frac{\overline{Y}}{\overline{x}} \ \overline{x}'\right]$

Where E_2 denotes the conditional expectation over second sample when the initial sample is fixed.

$$=E_1\left[\sum_{i=1}^{(n_n')} prob.\frac{\overline{Y}}{\overline{x}} \ \overline{x}'\right]$$

The probability of selecting the second sample for the fixed initial sample fixed is given by $\frac{\overline{x}}{\binom{n'}{\overline{x'}}}$

 $p_{i} = \frac{x_{i}}{X} = \frac{x_{i}}{\Sigma x_{i}} \text{ probability of choosing n from n'} (n'_{n})$ So. $x^{-1} = \frac{1}{(n'_{n})} \sum_{I}^{(n'_{n})} \overline{x}$ $E[\hat{Y}] = E_{1} \left[\sum_{I}^{\binom{n'}{n}} \frac{\overline{x}}{\binom{n'}{n}(x^{-1})} \frac{\overline{y}}{x} x^{-1} \right]$ $= E_{1} \left[\frac{I}{\binom{n'}{n}} \sum_{I}^{\binom{n'}{n}} \overline{y} \right]$ $= E_{1} \left[Y^{-1} \right]$ $= \overline{Y}$

Next
$$Var\left(\hat{\overline{Y}}\right) = E\left[\hat{\overline{Y}}^{2}\right] - \left(E\left[\hat{\overline{Y}}\right]\right)^{2}$$
$$E\left[\hat{\overline{Y}}^{2}\right] = E\left[\frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}} \, \overline{x}'^{2}\right]$$
$$= E_{1}E_{2}\left[\frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}} \, \overline{x}'^{2}\right]$$
$$= E_{1}\left[\sum_{l}^{\binom{n'}{n}} prob. \frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}} \, \overline{x}'^{2}\right]$$
$$= E_{1}\left[\sum_{l}^{\binom{n'}{n}} \frac{\overline{x}}{\binom{n'}{n}} \frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}} \, \overline{x}'^{2}\right]$$
$$= \frac{1}{\binom{n'}{n}} E_{1}\left[x^{-1}\sum_{l}^{\binom{n}{n}} \frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}}\right]$$
$$= \frac{1}{\binom{n'}{n}} \left[\sum_{l}^{\binom{n'}{n}} \frac{1}{\binom{n'}{n}} \left\{x^{-1}\sum_{l}^{\binom{n}{n}} \frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}}\right\}\right]$$
$$= \frac{1}{\binom{N}{n'}\binom{n'}{n}} \left[\sum_{l}^{\binom{n'}{n'}} x^{-1}\left\{\sum_{l}^{\binom{n}{n}} \frac{\overline{\overline{Y}}^{2}}{\overline{x}^{2}}\right\}\right]$$

There fore variance becomes,

$$Var[\hat{\overline{Y}}] = \frac{1}{\binom{N}{n'}\binom{n'}{n}} \left[\sum_{j=1}^{\binom{N}{n'}} x^{-1} \left\{ \sum_{j=1}^{\binom{n'}{n}} \frac{\overline{Y}^2}{\overline{x}^2} \right\} \right] - \overline{Y}^2$$

** Non-Sampling Errors:-

The theory of sampling scheme is assume that,

- (i) Some probability sampling scheme is used.
- (ii) The observation on the i^{th} unit of the population i.e. Y_i is the correct value.

Then the error in the estimate is surely due to random sampling this error is known as sampling error. In general, these are other type of errors, these errors are due to measurement feabalation, editing, etc. These errors are known as non-sampling errors. The main sources of non-sampling errors are,

- (i) Lake of precision in reporting observation.
- (ii) Incomplete coverage of the sample.
- (iii) Faulty method of estimation is used.

When a complete count is made there is no sampling error but there will be non-sampling errors.

As the sample size increases the sampling error will tend to decrease and non-sampling errors will tend to increase.

When a sample is selected both the type of errors will remain present.